1    Chromosome-scale *Salvia hispanica* L. (Chia) genome assembly reveals rampant *Salvia*

2    interspecies introgression

3    Julia Brose[1], John P. Hamilton[2,3], Nicholas Schlecht[4], Dongyan Zhao[1], Paulina M. Mejía-Ponce[5],

4    Arely Cruz Pérez[5], Brieanne Vaillancourt[2], Joshua C. Wood[2], Patrick P. Edger[6], Salvador Montes-

5    Hernandez[7], Guillermo Orozco de Rosas[8], Björn Hamberger[4], Angélica Cibrian Jaramillo[5], and C.

6    Robin Buell[2,3,9]

7    [1]Department of Plant Biology, Michigan State University, East Lansing, MI, USA

8    [2]Center for Applied Genetic Technologies, University of Georgia, Athens, GA, USA

9    [3]Department of Crop and Soil Sciences, University of Georgia, Athens, GA, USA

10    [4]Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing,

11    MI, USA

12    [5]National Laboratory for Genomics of Biodiversity (UGA-Langebio), CINVESTAV, Irapuato,

13    Guanajuato, Mexico

14    [6]Department of Horticulture, Michigan State University, East Lansing, MI, USA

15    [7]Campo Experimental Bajío, Instituto Nacional de Investigaciones Forestales, Agrícolas y

16    Pecuarias, Km 6.5 carretera Celaya-San Miguel de Allende, C.P. 38110, Celaya, Guanajuato,

17    México.

18    [8]CHIABLANCA SC DE RL, Acatic, Jalisco. Mexico

19    [9]Institute of Plant Breeding, Genetics, and Genomics, University of Georgia, Athens, GA, USA

20

21

22

23

**ABSTRACT**

*Salvia hispanica* L. (Chia), a member of the Lamiaceae, is an economically important crop in Mesoamerica, with health benefits associated with its seed fatty acid composition. Chia varieties are distinguished based on seed color including mixed white and black (Chia pinta) and black (Chia negra). To facilitate research on Chia and expand on comparative analyses within the Lamiaceae, we generated a chromosome-scale assembly of a Chia pinta accession and performed comparative genome analyses with a previously published Chia negra genome assembly. The Chia pinta and negra genome sequences were highly similar as shown by a limited number of single nucleotide polymorphisms and extensive shared orthologous gene membership. There is an enrichment of terpene synthases in the Chia pinta genome relative to the Chia negra genome. We sequenced and analyzed the genomes of 20 Chia accessions with differing seed color and geographic origin revealing population structure within *S. hispanica* and interspecific introgressions of *Salvia* species. As the genus *Salvia* is polyphyletic, its evolutionary history remains unclear. Using large-scale synteny analysis within the Lamiaceae and orthologous group membership, we resolved the phylogeny of *Salvia* species. This study and its collective resources further our understanding of genomic diversity in this food crop and the extent of inter-species hybridizations in *Salvia*.

**PLAIN LANGUAGE SUMMARY**

Chia pinta is an economically important crop due to the high fatty acid present in the seeds. There are multiple types of Chia based on the seeds color including mixed which and black (Chia pinta), black (Chia negra), and white (Chia blanca). We generated a genome assembly of Chia pinta and compared it to existing genome assemblies. While the assemblies are highly similar there are key differences in terpene synthase composition between Chia pinta and Chia negra. We also sequenced 20 other Chia accessions with different seed color and geographic origin to determine a population structure within Chia. We generated genomic resources to further our understanding of this food crop.

**ABBREVIATIONS**

BGC Biosynthetic gene cluster

BUSCO Benchmarking Universal Single Copy Orthologs

GO Gene ontology

SNP Single nucleotide polymorphism

TIR Terminal inverted repeat

56   TPS Terpene synthase

57   WGS Whole genome shotgun

58

59

## INTRODUCTION

61       Chia (*Salvia hispanica* L.) belongs to the largest genus within the Lamiaceae containing
62   approximately 980 species (Hu et al., 2018). Chia is a notable and economically important
63   species within the *Salvia* genus attributable to the high nutritional value of its seeds which
64   contain 16-26% protein, 23-41% fiber, and 20-34% polyunsaturated fatty acids, of which, 60% is
65   $\alpha$-lineolic acid (Muñoz et al., 2013). Historically, Chia was the third most economically
66   important crop in Mesoamerica, only behind maize and amaranth, due to its use in religious
67   practices and as a medicine (Valdivia-López & Tecante, 2015). The medicinal properties of Chia
68   include treatments for gastrointestinal, respiratory, urinary, obstetrics, skin, central nervous,
69   and ophthalmologic issues (Cahill, 2003). The traditional uses of Chia revolve around religious
70   practices which contributed to the decrease of Chia prominence and cultivation in the 15th
71   century following the invasion by conquistadors (Cahill, 2003). Chia was introduced to Spain
72   where it was named by Linnaeus as *Salvia hispanica* referencing the presumed origin of Spain
73   (Baldivia, 2018). While Chia originated in present day Mexico and Guatemala, it has since been
74   distributed throughout the world resulting in the emergence of diverse varieties (Cahill, 2004).

75       Chia varieties are characterized by their seed color and origin. The widely cultivated Chia
76   blanca has a white seed coat while Chia negra has a black seed coat that can occur in wild and
77   cultivated populations. Other seed coat colors include mixes of black and white seeds.
78   Morphological characteristics distinguishing cultivated from wild accessions mirror traits
79   observed in other domesticated species, such as decreased apical dominance, increased
80   branching, increased seed size, decreased pubescence, increased florescence length
81   determinism, increased anthocyanin pigmentation, variation in seed coat color and patterns,
82   increased plant height, and closed calyxes (Cahill, 2004). While phenotypically distinct, dietary
83   proteins are similar in wild and cultivated Chia accessions although wild accessions with higher
84   levels of polyunsaturated fatty acids have been reported (Peláez et al., 2019).

85       Robust genomic resources for the Lamiaceae facilitate comparative genomic analysis.
86   Within the Lamiaceae there are seven subfamilies with chromosome-scale genomes
87   [Ajugoideae, Callicarpoideae, Nepetoideae, Lamiodeae, Scutellariodeae, and Tectonoideae]
88   (Dong et al., 2018; Zhao et al., 2019a; b; Hamilton et al., 2020; He et al., 2022; Li et al., 2022;
89   Shen et al., 2022; Sun et al., 2022; Pan et al., 2023). Current genomic resources for Chia include
90   a genome assembly derived from an Australian black seeded variety (Chia negra; Wang et al.,
91   2022), a white seeded variety (Chia blanca; Li et al., 2023), and a Mexican Chia (Alejo-Jacuinde
92   et al., 2023) as well as transcriptomes constructed from wild and cultivated seeds (Peláez et al.,
93   2019). Expanding the number and diversity of chia accessions with genome assemblies and
94   sequence will facilitate our understanding of genetic diversity of this important crop as well as

4

95      provide resources for more informed breeding programs. In addition to diversity within Chia,

96      three other *Salvia* species occur in the same region in Mesoamerica (*Salvia uruapan* Fern.,

97      *Salvia tiliifolia* Vahl., and *Salvia polystachya* Ort.) that have similar uses as *S. hispanica* (Cahill,

98      2003). These species are challenging to distinguish from each other, but no reports indicate

99      hybridization with *S. hispanica*. A phylogeny of *Salvia,* based on 91 nuclear genes, places Chia

100     within *Salvia* sect. *Potiles* in a monophyletic clade (Lara-Cabrera et al., 2021). However, the

101     *Salvia* genus has yet to be fully resolved and remains polyphyletic with *S. tiliifolia* being placed

102     within two separate clades: the Angulatae and Polystachyae (Lara-Cabrera et al., 2021).

103     Therefore, additional phylogenetic analyses are necessary to achieve a comprehensive

104     resolution of the *Salvia* genus*.*

105        In this study, we report on the genome sequence of a Chia pinta accession, comparative

106     analyses with published Chia genomes, and analysis of genetic diversity in a set of 20 Chia

107     accessions revealing population structure between domesticated and wild Chia species and

108     evidence of interspecies hybridization of *S. tiliifolia* with Chia.

109     **RESULTS AND DISCUSSION**

110     **Chia Genome**

111        We selected a Chia pinta accession from Acatic**,** Jalisco, Mexico that produces mixed

112     color seeds and is grown as a superfood source throughout Mexico. Using 5.7 million PacBio

113     long reads (36.5 Gb) representing ~100x coverage of the predicted ~355 Mbp Chia genome

114     (Wang et al., 2022), we assembled the Chia pinta (2n=2x=12) genome using Canu (Koren et al.,

115     2017). Whole genome shotgun (WGS) reads were used to generate a k-mer (k=21) distribution

116     profile using GenomeScope indicating an estimated genome size of 338 Mbp with 62.6% unique

117     kmers and 0.5% heterozygosity. The initial Canu assembly was error corrected using the raw

118     PacBio reads using Arrow (Pacific Biosciences) followed by three rounds of error correction with

119     the Illumina WGS reads using Pilon (Walker et al., 2014). The error-corrected assembly

120     consisted of 2,094 contigs with a total length of 425.14 Mbp, which is substantially larger than

121     the previously estimated genome size. Haplotigs were removed from the assembly using

122     purgeHaplotigs (Roach et al., 2018) (-a = 50%) with an output consisting of "primary contigs"

123     representing the putative haploid genome sequence, "haplotigs" containing diverged

124     haplotypes, and "artefacts" representing contigs with very low or extremely high read

125     coverage. Following removal of haplotigs, the "primary contigs" size decreased from 425 Mbp

126     to 343 Mbp (Table 1). Manual examination of Chia vs. Chia self-alignments of contigs in the

127     'purged assembly' revealed five pairs of contigs that were putative residual haplotigs. Removal

128     of these contigs resulted in a 'purged assembly' containing 407 contigs with an N50 contig

129     length of 1.5 Mbp and a total size of 343.2 Mbp. The distribution of k-mers from WGS reads in

130    the final assembly was examined using KAT (Mapleson et al., 2017) revealing a single peak

131    indicating a haploid assembly with few retained haplotigs.

132         Using Hi-C sequence data, the contigs were assembled into six pseudomolecules,

133    consistent with the known chromosome number of Chia and the Chia negra Australian Black

134    (hereafter Chia negra) genome assembly (Wang et al., 2022). The final Chia pinta genome

135    assembly was 342 Mb final with an N50 of 62Mb, of which, 99.64% of the assembly was

136    anchored to one of the six pseudochromosomes (Table 1). Metrics for the final chromosome

137    assembly were calculated using only the six chromosomes. The GC content of the final assembly

138    was 36.6% consistent with the previously published Chia negra genome (Wang et al., 2022).

139    Alignment of Illumina WGS reads to the final assembly revealed 98.4% of the reads aligned to

140    the genome, of which, 99.5% were properly paired. Alignment of RNA-seq reads from a diverse

141    set of tissue types (leaf, inflorescence, stem, and root) showed an overall alignment rate

142    between 93.7% and 96.0%. To confirm the quality of the Chia pinta assembly, we used

143    Benchmarking Universal Single Copy Orthologs (BUSCO) (Simão et al., 2015) to determine the

144    representation of conserved orthologs in the final assembly. In total, 97.4% of the BUSCO

145    orthologs were complete with 86.6% as single copy, 10.8% duplicated, 0.7% fragmented, and

146    1.9% missing. Overall, these results indicate a high-quality Chia pinta genome assembly.

147    **Repetitive Sequences and Transposable Element Annotation in the Chia pinta genome**

148    Using *de novo* repetitive sequence identification with RepeatModeler coupled with sequences

149    from the Viridiplantae RepBase, RepeatMasker masked 46.8% of the Chia pinta genome. With

150    respect to transposable elements, retroelements were the dominant sequence with 40,151

151    retroelements occupying 15.15% of the Chia pinta genome while DNA transposons (36,807

152    elements) accounted for 4.86%. Unclassified interspersed repeats represented the largest

153    number of elements with 378,795 or 26.11% of the genome. The remaining repetitive elements

154    included rolling circle, small RNA, satellites, simple repeats, and low complexity sequences

155    make up less than 1% of the genome.

156         The Extensive *de-novo* TE Annotator (EDTA) was used to annotate the Chia pinta

157    genome for transposable elements revealing 314,306 elements spanning 149,780,410 bp

158    (43.64%) of the Chia pinta genome. Long terminal repeats comprise 21.33% of the genome, of

159    which, 5.7% were *Copia* elements and 11.45% were *Gypsy* elements; unknown long terminal

160    repeats comprise 4.13% of the genome. Terminal inverted repeat (TIR) sequences represent

161    20.01% of the genome with the largest portion (12.09%) belonging to Tc1_Mariner family. The

162    remaining TIRs are PIF_Harbinger (3.26%), hAT (2.32%), Mutator (1.80%), and CACTA (0.54%).

163    Helitrons are non-terminal inverted repetitive elements and comprise 2.3% of the genome.

164

**Annotation of the Chia Pinta Genome**

165

166    We annotated the Chia pinta genome for protein-coding genes resulting in 59,062

167    working gene models corresponding to 41,279 loci (Table 2). Working gene models had an

168    average transcript length of 3.1 kbp, coding sequence (CDS) length of 1,217 bp, exon length of

169    279 bp and intron length of 240 bp. Working gene models exhibited an average of 5.8 exons,

170    with 13.6% of transcripts being single-exon genes. The high confidence model set, a subset of

171    the working set which have expression and/or protein evidence, contains 53,053 gene models

172    representing 35,480 loci (Table 2). The high confidence set has an average transcript length of

173    3.3 kbp, exon length of 226 bp, intron length of 244 bp, and 6.1 exons per model; 6,105 gene

174    models are single exon models. We selected the longest model as a representative for each

175    gene locus from the working and high confidence model sets. With respect to BUSCO

176    representation, the high confidence representative models are 94.8% complete, of which,

177    84.8% are complete and single copy while 10% are complete and duplicated; 1.9% are

178    fragmented and 3.3% are missing. For the working representative models, 95.7% are complete

179    with 85.5% complete and single copy and 10.2% complete and duplicated; 1.7% fragmented

180    and 2.6% missing. Overall, the BUSCO results indicate a robust annotation of the Chia pinta

181    genome.

**Comparative Analyses of Chia Genome Assemblies**

182

183    There are currently three published long-read, chromosome-scale Chia genome

184    assemblies: Chia blanca (Li et al., 2023), Chia negra (Wang et al., 2022), and Mexican Chia

185    (Alejo-Jacuinde et al., 2023). BUSCO analysis of all three published Chia genomes revealed that

186    all of these assemblies were high quality and with robust gene annotation datasets. Syntenic

187    orthologs (syntelogs) were identified between all four assemblies revealing a high degree of

188    synteny between these genome assemblies (Figure 1) with limited disruptions that may be due

189    to assembly artifacts in the various genome assemblies. Due to the high degree of similarity

190    between the four Chia genomes, we performed detailed comparisons of our Chia pinta genome

191    to the chromosome-scale black seeded Chia negra in which 73.62% of the genes were colinear

192    within 1,178 syntenic blocks (Figure 1). Chia negra is a 344Mb genome assembly with 99.05%

193    anchored on to chromosomes and 3.3Mb unanchored (Wang et al., 2022) with 428 gaps,

194    amounting to a total of 191.2 kbp Ns. A total of 1,278,367 Single Nucleotide Polymorphisms

195    (SNPs) were identified between the Chia negra and Chia pinta genomes that were distributed

196    throughout the genome with 10.0% (127,210) residing in genic regions, 75.6% (967,385) in

197    intergenic regions, and 14.4% (184,772) within intronic regions of the Chia pinta genome.

198    Using Orthofinder with the predicted proteomes of both Chia pinta and Chia negra, we

199    identified 20,580 orthogroups, of which, 358 orthogroups (2,738 genes) were unique to Chia

200    pinta while 462 orthogroups (1,458 genes) were unique to Chia negra. Gene ontology (GO)

7

201 enrichment of the genes unique to Chia pinta revealed differences in certain biological process,
202 cellular components, and molecular function ontologies. Of particular interest was the
203 enrichment of the GO terms "defense response", and "diterpenoid biosynthetic process" with
204 45 terpene synthases identified in the GO terms "diterpenoid biosynthetic process" and
205 "terpene synthase activity".

206 BLASTP was used to search all representative proteins in Chia pinta and Chia negra
207 against a collection of known terpene synthases (TPSs). TPSs greater than 350 amino acids were
208 used to create a phylogeny to determine the relationships among the TPSs. After filtering, a
209 total of 111 TPSs in Chia pinta and 53 in Chia negra were identified. To confirm that this is not
210 due to annotation errors, Chia pinta TPS transcript sequences were used in a BLASTN search
211 against the Chia negra genome; no additional terpene synthases were identified in Chia negra
212 indicating these sequences are absent in the Chia negra genome assembly. A phylogeny was
213 constructed with putative TPS protein sequences from Chia pinta, Chia negra, and functionally
214 characterized TPSs to assign Chia TPSs to closest known functionally characterized TPSs. Despite
215 GO enrichment annotation of 'diterpenoid biosynthetic process', most enriched TPSs are within
216 the TPS-a and to a lesser degree TPS-b subfamilies which produce sesqui- and monoterpenes,
217 indicating an expansion of volatile terpenes. The discrepancy on the GO terms claiming
218 diterpenoid processes yet finding sesqui- and monoterpene synthases can be explained by GO
219 enrichment often misannotated TPSs as diTPSs.

220 The TPS-a subfamily contains 56 putative TPSs in Chia pinta and only four in Chia negra.
221 Of the 56 putative Chia pinta TPSs, 38 were found to enriched relative to Chia negra. The
222 enriched TPSs reside in clades that do not contain a Chia negra TPS. To further understand the
223 genomic context of the enriched TPSs, biosynthetic gene clusters (BGCs) membership and
224 synteny were used. There are 16 BGCs containing TPSs in Chia pinta present on chromosomes
225 1, 2, 3, 4, and 6. Notably, six of these BGCs contain 23 out of the 56 Chia pinta specific TPS-a
226 subfamily genes (Figure 2). This coincides with the expansion of the TPS-a subfamily in Chia
227 pinta. All Chia pinta enriched TPS-a BGCs contain syntenic genes between Chia pinta, Chia
228 negra, and *S. miltiorrhiza* (Figure 2). However, Chia pinta only shares one syntenic TPS with Chia
229 negra and three syntenic TPSs with *S. miltiorrhiza.* Many of the TPSs present in Chia pinta's
230 BGCs appear to be tandem duplications, most notably in the teal and green BGCs (Figure 2)**.**
231 However, some of the TPSs present in the green BGC are less than 350 amino acids indicating
232 they may be truncated.

233 The origin and expansions of TPS-a genes were examined through synteny with *S.*
234 *miltiorrhiza*. Two separate BGCs, purple and orange, contain paralogous TPSs yet are in distinct
235 syntenic blocks (Figure 2). Work in *S. miltiorrhiza* characterized orthologs of these genes (89%
236 identity) as (-)-5-epi-eremophilene synthases in which three TPSs (*SmSTPS1*, *SmSTPS2*, and

237   *SmSTPS3*) had differential gene expression yet identical biochemical activity (Fang et al., 2017).
238   The purple BGC contains one TPS that is a syntelog of *SmSTPS1*, but there are no syntelogs of
239   *SmSTPS2* or *SmSTPS3* (Figure 2) suggesting that a single gene was maintained and was tandemly
240   duplicated or that structural rearrangements occurred disrupting synteny with *SmSTPS2* or
241   *SmSTPS3*. The orange BGC contains TPSs that are equally related to *SmSTPS1* but are not
242   syntenic with the *S. miltiorrhiza* SmSTPS cluster. Instead, the homologs have moved into a
243   different syntenic block entirely. Additionally, there is a notable difference in gene expression
244   profiles of the purple and orange BGCs with the orange BGC largely expressed in the leaf and
245   stem whereas the purple clade has its highest expression in roots amongst the different
246   paralogs (Figure 2). This may exemplify how a BGC can evolve by duplication and
247   subfunctionalization resulting in distinct spatial gene expression patterns. The teal and yellow
248   BGCs indicate that there are no syntenic TPSs in *S. miltiorrhiza*. The minor enrichment in TPS-b
249   genes present in Chia pinta is largely due to expansion of a single clade. The closest functionally
250   characterized enzyme to this expanded clade was and (−)-exo-α-bergamotene synthase, having
251   between 62-67% identity for this clade.

252        Finding such a large difference in TPS-a abundance and identifying many of them within
253   BGCs between Chia pinta and Chia negra further supports the diversity that exists not just
254   within the *Salvia* genus, but even within Chia accessions. One potential source of the TPS
255   expansion could be due to sequencing gaps in the Chia negra genome assembly. Specially, there
256   are gaps in the purple BGC region of the Chia negra genome sequence. Therefore, these TPSs
257   could be present within the species, but were not captured by the genome assembly. However,
258   for the remaining five BGCs there are no assembly gaps in the Chia negra genome assembly and
259   when the predicted transcripts for the TPSs were searched against the Chia negra genome,
260   there were no hits for these regions. To determine if the TPSs are unique to Chia pinta, we
261   examined the BGCs for syntelogs in the two other long-read Chia genome assemblies. The teal,
262   orange, pink, green, and yellow BGCs contain syntelogs in Chia pinta, Chia blanca, and Mexican
263   Chia whereas the purple BGC contains only syntelogs between Chia pinta and Mexican Chia.
264   Thus, diversity in TPSs is present between Chia accessions suggesting variation in terpenoid
265   profiles that may be associated with local adaptation.

**Lamiaceae Phylogeny and Gene Family Expansions**

267        To determine the evolutionary relationships of Lamiaceae species with Chia pinta, a
268   species phylogeny was constructed using high-quality available genome sequences from 23
269   species from seven tribes in the Lamiaceae (Figure 3). Using the multiple sequence alignment
270   option in Orthofinder, 923,746 genes were assigned to orthogroups. As shown in Figure 3, the
271   Nepetoideae tribe is sister to Ajugoideae, Lamiodeae, and Scutellariodeae, the Callicarpoideae
272   and Tectonoideae are sister to all other species, and the Premnoideae is sister to all other

273  subfamilies. The relationships between the tribes in this genome-derived tree differs from a
274  published phylogeny derived from 520 single copy transcripts (Godden et al., 2019) in which the
275  Nepetoideae is sister to Ajugoideae, Lamiodeae, Scutellariodeae, Premnoideae, and
276  Tectonoideae. The topology difference between these two phylogenetic estimates could be due
277  to a combination of species sampling and data quality differences.

278      Gene expansions and contractions of single copy orthologs throughout the Lamiaceae
279  were identified using CAFE (Figure 3A) and placed on the species tree phylogeny revealing large
280  expansions and contractions throughout the Lamiaceae. The node branching of the
281  Nepetoideae indicates a gene family expansion of 1,506 genes and contraction of 1,688 genes.
282  The branch point from *S. hispanica* and *Salvia splendens* reveals 2,901 gene expansions and
283  12,295 gene contractions indicating substantial differences within the *Salvia* genus.

284      Synteny between genomes serve as a tool for examining evolution reflecting ancestral
285  conservation of gene order. Using Chia pinta as the reference genome, we examined synteny
286  within 11 chromosome-scale assemblies, spanning six tribes of the Lamiaceae family, revealing
287  extensive conservation among the genomes (Figure 3B). In total, 182 Chia pinta genes were
288  found to have a one-to-one syntenic relationship across all 11 species.

289      The polyphyletic nature of *Salvia* is highlighted by orthogroup membership. Of the
290  39,379 orthogroups containing 211,888 genes there were 12,987 orthogroups, containing
291  165,520 genes, in common among all *Salvia* (Figure 4A). The next highest number of
292  orthogroups are unique to *S. rosmarinus* closely followed by *S. officinalis* and then *S. splendens*
293  (Figure 4A). We also performed syntenic analyses between the genomes of four *Salvia* species
294  to further our understanding of the species relationship in this polyphyletic genus. As expected,
295  Chia pinta shares extensive synteny with other *Salvia* species (Figure 4b). *S. splendens* is
296  reported to be a tetraploid (Jia et al., 2021). Based on orthogroup membership, 25% (4,684) of
297  orthogroups shared by *S. splendens* and Chia pinta contain two *S. splendens* genes for each Chia
298  pinta gene. This pattern reflects that *S. splendens* is a tetraploid and Chia pinta is a diploid.
299  There are also two syntenic blocks in *S. splendens* for each block within Chia pinta, the syntenic
300  blocks exist across four chromosomes in *S. splendens* (Figure 4b and 4c)*. It has been reported
301  that there is a single shared whole genome duplication between Chia pinta and *S. splendens*
302  and an additional duplication just in *S. splendens* (Jia et al., 2021; Wang et al., 2022). Therefore,
303  the four unique chromosomes syntenic to a single chromosome in Chia pinta could be due to
304  chromosomal fusions in Chia pinta or chromosomal fissions in *S. splendens.* Within the *Salvia*
305  genus there are large regions of fragmented synteny between Chia pinta and *S. officinalis* as
306  well as between *S. splendens* and *S. officinalis.* The fragmentation could be present due to
307  different ancestry of Chia pinta and *S. officinalis*. As *Salvia* is a polyphyletic genus (Lara-Cabrera
308  et al., 2021), this could be indicative of how distantly related these two species are. An

309 alternative hypothesis is that they share a common ancestor, but the divergence time between
310 species is so long that conserved genetic regions have been differentially fractionated (i.e.
311 unique gene loss patterns). This is consistent with the large gene family expansions and
312 contractions in the node that splits the *Salvia* species*.*

**Population Structure of Chia**

314 Seed coat color is a frequent descriptor for Chia accessions with Chia white seeded
315 blanca varieties while Chia negra, Chia cualac, and Chia xonostli are predominately black-
316 seeded (Figure 5a). Chia pinta seeds are a mix of both black and white seeds (Figure 5a). A
317 diversity panel of 19 Chia accessions including wild and cultivated accessions along with two *S.*
318 *tiliifolia* accessions with origins throughout Mexico was constructed and sequenced to reveal
319 genetic diversity among accessions and provide insight into population structure of cultivated
320 and wild Chia varieties. The percentage of reads aligned to the Chia pinta genome ranged from
321 95.5% to 97.7% for the *S. tiliifolia* samples and 96.3%-98.5% for the Chia varieties suggesting
322 that the two species share substantial sequence similarity. Population structures were inferred
323 with admixture with k=2 to k =13. Population structure admixture results suggested through
324 the cross-validation error plot that there are two possible number of populations: four and nine
325 as the local minima being at four and the global minima at nine in the cross-validation error
326 plot.

327 Using a k=4, broad population groups are present that can be assigned to known
328 categories of Chia: Chia pinta (yellow), *S. tiliifolia* (purple), Chia negra and Chia Xonostli (blue),
329 and Chia Cualac (pink). The population structure indicates that the phenotypic and origin
330 grouping reflects the genetic structure of the population. Chia pinta accessions are
331 domesticated Chia varieties whereas Chia negra and Chia Xonostli are classified as wild due to
332 their more open calyx and other wild traits. Chia negra is in the same population group with the
333 less widely known Chia Xonostli which is similar to Chia negra yet categorized differently due to
334 its domesticated traits. Historically, Chia Xonostli was found in the states of Jalisco, Guanajuato,
335 Veracruz, and Hidalgo. Chia Cualac was reported to be semi-domesticated and forms their own
336 group with some admixture from Chia Xonostli (Peláez et al., 2019). This follows the hypothesis
337 that wild introgressions are present throughout the populations. One *S. tiliifolia* accession is
338 admixed with Chia pinta. *S. tiliifolia* is nearly indistinguishable from Chia and is known to grow
339 in the same areas as Chia pinta; thus, it is possible that these species hybridize and form a
340 population of *S. tiliifolia* that is highly admixed with Chia pinta. Feral hybrid accessions could
341 continue to evolve through hybridization with domesticated Chia yielding the admixture
342 present within one accession of *S. tiliifolia* (Figure 5)*.*

343

11

**CONCLUSIONS**

In this study, a high-quality chromosome-scale genome assembly of Chia pinta was generated that allowed for additional genomic comparisons within the economically important crop including three other long-read, chromosome-scale Chia assemblies that showed extensive synteny among the genome sequences. Comparative genomic tools were used to determine differences within Chia accessions and throughout the Lamiaceae. Interestingly, Chia pinta was enriched in TPSs and contains novel TPSs compared to the Chia negra with some TPSs located within BGCs and syntenic with *S. miltiorrhiza*. Further examination of TPSs within BGCs among the four Chia genome assemblies revealed further diversification suggestive of variation in terpenoid biosynthesis among varieties.  Through sequencing of a diversity panel, the population structure of Chia revealed introgression with other *Salvia* species.

**MATERIALS AND METHODS**

**Plant materials**

Different Chia varieties were collected throughout Mexico. Plants were grown in an experimental field in Celaya, Guanajuato, Mexico (20.578°, −100.822°) at the Instituto Nacional de Investigaciones Forestales, Agrícolas y Pecuarias (INIFAP).

**Nucleic acid isolation, library construction, and sequencing**

For construction of a reference genome, DNA was isolated from medium-sized leaves from a mature plant (13.5 weeks old) of accession SM_ACJ2017 using a modified protocol from Doyle and Doyle (1987) and Healey et al. (2014). Large insert (>15kb, >20kb) PacBio libraries were made with the SMRTbell™ Template Prep Kit and sequenced on the PacBio Sequel platform at the University of Georgia, Georgia Genomics and Bioinformatics Core (GGBC, UG Athens, GA, RRID:SCR_010994). A whole genome shotgun library for reference error correction was prepared using the Illumina TruSeq Nano DNA Library Preparation Kit and sequenced in paired-end mode, 150 nt in length on a HiSeq 4000 at the Michigan State University Research Technology Support Facility (RTSF). Whole genome shotgun libraries for use in error correction and diversity panel variant analyses were constructed as described previously in Hardigan *et al.* 2016 (Hardigan et al., 2016) and sequenced at the Michigan State University RTSF in paired-end mode on a HiSeq4000 generating 150 nt reads. RNA was isolated from three biological replicates from a core set of tissues (leaf, inflorescence, lateral stem, secondary root) from the reference accession SM_ACJ2017 as described previously in Peláez *et al.* 2019 (Peláez et al., 2019). RNA-seq libraries were prepared using the Illumina TruSeq Stranded mRNA Library Preparation Kit and sequenced on an Illumina HiSeq 4000 generating 150 nt paired end reads for one replicate and 50 nt single end reads for the other two replicates; library preparation and sequencing were performed at the Michigan State University Research Technology Support

379  Facility (RTSF). A Phase Genomics Proximo Hi-C library was prepared from Chia pinta leaf tissue

380  and sequenced by Phase Genomics (Seattle, WA) on the NextSeq 500 generating paired end

381  150 nt reads.

**Chia pinta genome assembly**

383  PacBio reads greater than 10 kbp (1.2 million reads, 21.6 Gb) were used to generate the initial

384  assembly using Canu (v1.7; Koren et al., 2017) with a corrected ErrorRate of 0.15%. The initial

385  assembly was polished with the raw PacBio reads using Arrow in the SMRT Analysis package

386  (v5.0.1.9585; Pacfici Biosciences), followed by three rounds of error correction with 56 million

387  Illumina WGS reads (150 nt paired-end WGS reads, 45X coverage) using Pilon (v1.22; Walker et

388  al., 2014). Potential haplotigs were purged using purgeHaplotigs (v1.0.4; Roach et al., 2018)

389  with the "maximum match score (-m)" of 500% and "-a = 50% ". Contigs were scaffolded to a

390  chromosome scale assembly using Hi-C reads and Proximo pipeline with an input chromosome

391  number of six by Phase Genomics (Bickhart et al., 2017). Scaffolded contigs were visualized with

392  Juicebox (v1.9.8; Durand et al., 2016).

**Genome annotation**

394  A custom repeat library (CRL) was generated using RepeatModeler (v2.0.1;Flynn et al.,

395  2020) and protein coding genes were removed from the CRL using ProteinExcluder (v1.2;

396  Campbell et al., 2014). The Viridiplantae RepBase repeats (v20150807) were then added to

397  create the final CRL. The genome assembly was hard and soft masked using RepeatMasker

398  (v4.1.0; Smit et al.) with the CRL with the parameters: -s -nolow -no_is. RNA-seq libraries were

399  cleaned using Cutadapt (v2.9; Martin, 2011) (--times 2 --minimum-length 100 --quality-cutoff

400  10) and then aligned to the genome assembly with HISAT2 (v2.2.0;Kim et al., 2019) (--max-

401  intronlen 5000 --rna-strandness RF –dta –no-unal). The RNA-seq alignments were then

402  assembled into transcript assemblies using Stringtie (v2.1.1; Kovaka et al., 2019).

403  *Ab initio* gene models were predicted on the soft-masked genome assembly using the

404  BRAKER2 pipeline (v2.1.5; Brůna et al., 2021) using the leaf RNA-seq library CHI_AA as a source

405  for hints. The *ab initio* gene models were then refined using PASA2 (v2.4.1; Campbell et al.,

406  2006) with the RNA-seq transcript assemblies as a source of transcript evidence to produce the

407  working gene model set. High confidence gene models were selected from the working gene

408  model set by first calculating working gene model abundances of the RNA-seq libraries for the

409  working gene models with Kallisto (v0.46.0; Bray et al., 2016), then searching the working gene

410  models against PFAM (v32.0; Mistry et al., 2021) with HMMER (v3.2.1; Mistry et al., 2013).

411  Working gene models with a TPM >1 in at least one RNA-seq library or a non-transposable

412  element related PFAM domain match and no partial or containing an internal stop codon were

413  identified as high confidence gene models. Functional annotation was assigned to the working

414  gene model by searching the protein sequences against the Arabidopsis proteome (TAIR10),

415  PFAM (v32.0; Mistry et al., 2021) and the Swiss-Prot plant proteins (release 2015_08). Search

416  results were processed in the same order and the function of the first hit encountered was

417  assigned to the gene model. Repetitive elements were identified using EDTA (v2.1.0; Ou et al.,

418  2019) with the parameters species set to "others" and step set to "all".

**Genome quality assessment**

420  Quality assessment of the genome assembly was performed by aligning WGS reads cleaned for

421  low quality bases and adaptors using Cutadapt (v3.4; Martin, 2011) to the final assembly using

422  BWA-mem (v0.7.16a; Li, 2013). Assemblathon.pl

423  (https://github.com/KorfLab/Assemblathon/blob/master/assemblathon_stats.pl) was used to

424  generate genome metrics. BUSCO (v3.1.0.Py3; Simão et al., 2015) embryophyta_odb10 was

425  used to determine genic representation in the final assembly. Jellyfish (v.2.3.0; Marçais &

426  Kingsford, 2011) with the option -m 21 was used to count kmers that were then visualized in

427  GenomeScope (v2.0; Ranallo-Benavidez et al., 2020) with kmer length 21 was used to verify

428  genome size and heterozygosity from the WGS reads from Chia pinta (CHI_AN). The Kmer

429  Analysis Toolkit (v2.4.1; Mapleson et al., 2017) was used to examine the assembly for retained

430  haplotigs. Synteny between the chia genome assemblies (Wang et al., 2022; Alejo-Jacuinde et

431  al., 2023; Li et al., 2023) was analyzed using GENESPACE (v.1.1.10;Lovell et al., 2022). Syntenic

432  comparison between Chia pinta and Chia negra was also performed using MCScanX (Wang et

433  al., 2012).

**Lamiaceae phylogeny and comparative analysis**

435  Publicly available genomes of *Callicarpa armericana* (Hamilton et al., 2020), *Cleorodendrum*

436  *inerme* (He et al., 2022) , *Hyssopus officinalis* (Lichman et al., 2020), *Nepeta cataria* (Lichman et

437  al., 2020), *Nepeta mussinii* (Lichman et al., 2020), *Ocimum basilicum* (Bornowski et al., 2020),

438  *Origanum majorana* (Bornowski et al., 2020), *Origanum vulgare* (Bornowski et al., 2020), *Perilla*

439  *frustescens*(Zhang et al., 2021; Tamura et al., 2022), *Pogostemon cablin* (Shen et al., 2022),

440  *Salvia miltiorrhiza* (Pan et al., 2023), *Salvia officinalis* (Li et al., 2022), *Salvia rosmarinus*

441  (Bornowski et al., 2020), *Salvia splendens* (Jia et al., 2021), *Scutellaria baicalensis* (Zhao et al.,

442  2019b), *Scutellaria barbata* (Xu et al., 2020), *Tectona grandis* (Zhao et al., 2019a), *Thymus*

443  *quinquecostatus* (Sun et al., 2022)*, Lavandula angustifolia* (Hamilton et al., 2023) and *Premna*

444  *obstusifolia* (He et al., 2022) were obtained and quality assessed using BUSCO (v5.5.0; Simão et

445  al., 2015) embryophyta_odb10. Species with genome assembly complete BUSCO scores greater

446  than 90% and annotation complete BUSCO scores greater than 80% were used in further

447  comparative analysis. Orthogonal genes and species tree phylogeny were built using

448  OrthoFinder (v.2.5.4; Emms & Kelly, 2019) with options -M msa -T raxml. The species tree

449  output was covered into an ultrametric tree using the make_ultrametric command in

450   OrthoFinder (v.2.5.4; Emms & Kelly, 2019). Branch lengths were rescaled using the *Premna*
451   *obstusifolia* divergence date of 16.06 MYA retrieved from the TimeTree of Life resource (Kumar
452   et al., 2022). Gene family expansions and contractions were identified using CAFE (v.4.2.1; Han
453   et al., 2013) with the following scripts with default parameters: cafetutorial_report_analysis.py
454   and cafetutorial_draw_tree.py. Syntelogs through the Lamiaceae were obtained for the
455   chromosome scale assemblies within the Lamiaceae and visualized using GENESPACE (v.1.1.10;
456   Lovell et al., 2022).

**Gene ontology term enrichment**

458   Gene ontology (GO) terms were assigned to high confidence Chia pinta genes using
459   InterProScan (v5.63-95.0; Jones et al., 2014). GO descriptions were added using the
460   ontologyIndex package (Greene et al., 2017) and enrichment was calculated using the topGO R
461   package (Alexa & Rahnenfuhrer, 2010). GO terms with an FDR adjusted p-value < 0.05 were
462   considered significant.

**Terpene synthase identification**

464   BGCs were identified in Chia pinta, Chia negra, and *S. miltiorrhiza* with PlantiSMASH (Kautsar et
465   al., 2017). Enriched TPSs identified in the various BGCs were searched with NCBI BLAST the
466   nonredundant protein database to identify the closest functionally characterized TPSs. To
467   extract all TPSs from the genome, the high confidence representative protein models were
468   blasted against a reference set of known TPSs enzymes representing TPSs across all subfamilies.
469   The BLAST hits with an E-value 1E-5 or better were selected. These gene models were filtered
470   to remove any sequences smaller than 350 amino acids to ensure a quality phylogeny and
471   minimize pseudogenes. The final set of putative and reference TPS sequences were aligned
472   using clustal omega (v1.2.4; Sievers et al., 2011). A phylogenetic tree of the alignment was built
473   via RAXML (v8.2.12; Stamatakis, 2014) with the PROTGAMMA AUTO model, algorithm a, and
474   1000 bootstraps. Gene expression of terpene synthases was calculated using the single end
475   RNA-seq libraries and Cufflinks (v.2.2.1; Roberts et al., 2011) with the options -b and -u to
476   generate FPKM values for all Chia pinta genes. Orthologous genes from Chia pinta, Chia negra
477   (Wang et al., 2022) , Chia blanca (Li et al., 2023)and the Mexican Chia variety (Alejo-Jacuinde et
478   al., 2023) were identified using OrthoFinder (v.2.5.4; Emms & Kelly, 2019) with options -M msa -
479   T raxml.

**Population structure analysis**

481   Whole genome shotgun reads from the diversity panel were cleaned using Cutadapt (v3.4;
482   Martin, 2011) and aligned to the Chia genome using BWA-mem (v0.7.16a; Li, 2013). PicardTools
483   (v2.20.8; Picard toolkit, 2019) commands SortSam, MarkDuplicates, BuildBamIndex, and
484   CollectAlignmentSummaryMetrics were used to sort, convert files, and generate alignment

485  metrics. The GATK (v4.1.2.0; Van der Auwera & O'Connor, 2020) HaplotypeCaller with default
486  parameters was used to call variants. GenomicsDBImport with default parameters was used to
487  merge the varieties into a single VCF file and genotyped using GenotypeGVCFs.Separated. SNPs
488  were selected using the SelectVariants command. Hard filtering of the SNPs was performed
489  using the parameters QD < 2.0, QUAL < 30.0, SOR > 3.0, FS > 60.0, MQ < 40.0, MQRankSum < -
490  12.5, MQRankSum-12.5, ReadPosRankSum < -8.0. Additional filtering was performed using
491  VCFTools (v0.1.16; Danecek et al., 2011) with filtering –freq2 and –max-alleles 2 to retain only
492  bialleleic sites, minor allele frequency of 0.071, --max-missing 0.9, --minQ 30, --min-meanDP 15,
493  --max-meanDP 39.

494      SNPs were called relative to the Chia negra reference genome (Wang et al., 2022) using
495  nucmer from MUMmer (v4.0; Marçais et al., 2018) with the options –maxgap=2500, --
496  minmatch=11, and --mincluster=25. SNPs were quality filtered using the delta-filter command
497  in MUMmer with the -r flag. (v4.0; Marçais et al., 2018). The SNP set from the diversity panel
498  and from the alignment of the two genome assemblies were combined and converted into bed
499  format using PLINK 2.0 (v.alpha2.3; Purcell & Chang; Chang et al., 2015) resulting in 156,829
500  total SNPs. Population structure was inferred with Admixture (v.1.3.0; Alexander et al., 2009)
501  and a SNP phylogenetic tree built with SNPhylo (v.20160204; Lee et al., 2014) using default
502  parameters.

503  **ACKNOWLEDGMENTS**

507  **CONTRIBUTIONS**

508  ACJ and CRB conceived of the study. SM-HGOR, PMM-P and ACP collected samples. PMM-P and
509  ACP prepared materials. JB, JPH, NS, DZ, JCW and BV performed data analyses and drafted the
510  manuscript. PPE, BH, ACJ, and CRB supervised and performed project administration. All
511  authors approved of the manuscript.

512
513  **DATA AVAILABILITY STATEMENT**

514  The raw sequence reads are available in the National Center for Biotechnology Information
515  Sequence Read Archive under BioProject PRJNA744892. The genome assembly, annotation, and
516  large data sets (genome assembly, genome annotation) reported in this study are available in
517  Figshare via 10.6084/m9.figshare.24546049.

518  **CONFLICT OF INTERESTS**

519     The authors declare no conflict of interests.

520

## REFERENCES

522     Alejo-Jacuinde, G., Nájera-González, H.R., Chávez Montes, R.A., Gutierrez Reyes, C.D., Barragán-
523        Rosillo, A.C., Perez Sanchez, B., Mechref, Y., López-Arredondo, D., Yong-Villalobos, L., &
524        Herrera-Estrella, L. (2023). Multi-omic analyses reveal the unique properties of chia (Salvia
525        hispanica) seed metabolism. *Communications Biology*, *6*. https://doi.org/10.1038/s42003-
526        023-05192-4

527     Alexa, A., & Rahnenfuhrer, J. (2010). topGO: Enrichment Analysis for Gene Ontology

528     Alexander, D.H., Novembre, J., & Lange, K. (2009). Fast model-based estimation of ancestry in
529        unrelated individuals. *Genome Research*, *19*, 1655–1664.
530        https://doi.org/10.1101/gr.094052.109

531     Van der Auwera, G.A., & O'Connor, B. (2020). *Genomics in the Cloud: Using Docker, GATK, and
532        WDL in Terra*. 1st Editio. O'Reilly Media.

533     Baldivia, A.S. (2018). A Historical Review of the Scientific and Common Nomenclature Associated
534        with Chia: From *Salvia hispanica* to *Salvia mexicana* and Chian to Salba. *Agricultural
535        Research & Technology: Open Access Journal*, *18*.
536        https://doi.org/10.19080/artoaj.2018.18.556047

537     Bickhart, D.M., Rosen, B.D., Koren, S., Sayre, B.L., Hastie, A.R., Chan, S., Lee, J., Lam, E.T., Liachko,
538        I., Sullivan, S.T., Burton, J.N., Huson, H.J., Nystrom, J.C., Kelley, C.M., Hutchison, J.L., Zhou,
539        Y., Sun, J., Crisà, A., Ponce De León, F.A., Schwartz, J.C., Hammond, J.A., Waldbieser, G.C.,
540        Schroeder, S.G., Liu, G.E., Dunham, M.J., Shendure, J., Sonstegard, T.S., Phillippy, A.M., Van
541        Tassell, C.P., & Smith, T.P.L. (2017). Single-molecule sequencing and chromatin
542        conformation capture enable de novo reference assembly of the domestic goat genome.
543        *Nature Genetics*, *49*, 643–650. https://doi.org/10.1038/ng.3802

544     Bornowski, N., Hamilton, J.P., Liao, P., Wood, J.C., Dudareva, N., & Buell, C.R. (2020). Genome
545        sequencing of four culinary herbs reveals terpenoid genes underlying chemodiversity in the
546        Nepetoideae. *DNA Research*, *27*, 1–12. https://doi.org/10.1093/dnares/dsaa016

547     Bray, N.L., Pimentel, H., Melsted, P., & Pachter, L. (2016). Near-optimal probabilistic RNA-seq
548        quantification. *Nature Biotechnology*, *34*, 525–527. https://doi.org/10.1038/nbt.3519

549  Brůna, T., Hoff, K.J., Lomsadze, A., Stanke, M., & Borodovsky, M. (2021). BRAKER2: Automatic
550      eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein
551      database. *NAR Genomics and Bioinformatics*, *3*, 1–11.
552      https://doi.org/10.1093/nargab/lqaa108

553  Cahill, J.P. (2003). Ethnobotany of Chia, *Salvia hispanica* L. (Lamiaceae). *Economic Botany*, *57*,
554      604–618. https://doi.org/10.1663/0013-0001(2003)057[0604:EOCSHL]2.0.CO;2

555  Cahill, J.P. (2004). Genetic diversity among varieties of Chia (*Salvia hispanica* L.). *Genetic*
556      *Resources and Crop Evolution*, *51*, 773–781.
557      https://doi.org/10.1023/B:GRES.0000034583.20407.80

558  Campbell, M.A., Haas, B.J., Hamilton, J.P., Mount, S.M., & Robin, C.R. (2006). Comprehensive
559      analysis of alternative splicing in rice and comparative analyses with Arabidopsis. *BMC*
560      *Genomics*, *7*, 1–17. https://doi.org/10.1186/1471-2164-7-327

561  Campbell, M.S., Law, M.Y., Holt, C., Stein, J.C., Moghe, G.D., Hufnagel, D.E., Lei, J.,
562      Achawanantakun, R., Jiao, D., Lawrence, C.J., Ware, D., Shiu, S.H., Childs, K.L., Sun, Y., Jiang,
563      N., & Yandell, M. (2014). MAKER-P: A Tool kit for the rapid creation, management, and
564      quality control of plant genome annotations. *Plant Physiology*, *164*, 513–524.
565      https://doi.org/10.1104/pp.113.230144

566  Chang, C.C., Chow, C.C., Tellier, L.C.A.M., Vattikuti, S., Purcell, S.M., & Lee, J.J. (2015). Second-
567      generation PLINK: Rising to the challenge of larger and richer datasets. *GigaScience*, *4*, 1–
568      16. https://doi.org/10.1186/s13742-015-0047-8

569  Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., Handsaker, R.E.,
570      Lunter, G., Marth, G.T., Sherry, S.T., McVean, G., & Durbin, R. (2011). The variant call format
571      and VCFtools. *Bioinformatics*, *27*, 2156–2158.
572      https://doi.org/10.1093/bioinformatics/btr330

573  Dong, A.X., Xin, H.B., Li, Z.J., Liu, H., Sun, Y.Q., Nie, S., Zhao, Z.N., Cui, R.F., Zhang, R.G., Yun, Q.Z.,
574      Wang, X.N., Maghuly, F., Porth, I., Cong, R.C., & Mao, J.F. (2018). High-quality assembly of
575      the reference genome for scarlet sage, *Salvia splendens,* an economically important
576      ornamental plant. *GigaScience*, *7*, 1–10. https://doi.org/10.1093/gigascience/giy068

577  Doyle, J.J., & Doyle, J.L. (1987). A rapid DNA isolation procedure from small quantities of fresh
578      leaf tissues.. *Pytochemical Bulletin*, *19*, 11–15

579  Durand, N.C., Robinson, J.T., Shamim, M.S., Machol, I., Mesirov, J.P., Lander, E.S., & Aiden, E.L.
580      (2016). Juicebox Provides a Visualization System for Hi-C Contact Maps with Unlimited
581      Zoom. *Cell Systems*, *3*, 99–101. https://doi.org/10.1016/j.cels.2015.07.012

582     Emms, D.M., & Kelly, S. (2019). OrthoFinder: Phylogenetic orthology inference for comparative
583         genomics. *Genome Biology*, *20*, 238. https://doi.org/10.1186/s13059-019-1832-y

584     Fang, X., Li, C.Y., Yang, Y., Cui, M.Y., Chen, X.Y., & Yang, L. (2017). Identification of a novel (-)-5-
585         epieremophilene synthase from salvia miltiorrhiza via transcriptome mining. *Frontiers in*
586         *Plant Science*, *8*, 1–11. https://doi.org/10.3389/fpls.2017.00627

587     Flynn, J.M., Hubley, R., Goubert, C., Rosen, J., Clark, A.G., Feschotte, C., & Smit, A.F. (2020).
588         RepeatModeler2 for automated genomic discovery of transposable element families.
589         *Proceedings of the National Academy of Sciences of the United States of America*, *117*, 9451–
590         9457. https://doi.org/10.1073/pnas.1921046117

591     Godden, G.T., Kinser, T.J., Soltis, P.S., Soltis, D.E., & Chaw, S.M. (2019). Phylotranscriptomic
592         Analyses Reveal Asymmetrical Gene Duplication Dynamics and Signatures of Ancient
593         Polyploidy in Mints. *Genome Biology and Evolution*, *11*, 3393–3408.
594         https://doi.org/10.1093/GBE/EVZ239

595     Greene, D., Richardson, S., & Turro, E. (2017). OntologyX: A suite of R packages for working with
596         ontological data. *Bioinformatics*, *33*, 1104–1106.
597         https://doi.org/10.1093/bioinformatics/btw763

598     Hamilton, J.P., Godden, G.T., Lanier, E., Bhat, W.W., Kinser, T.J., Vaillancourt, B., Wang, H., Wood,
599         J.C., Jiang, J., Soltis, P.S., Soltis, D.E., Hamberger, B., & Robin Buell, C. (2020). Generation of
600         a chromosome-scale genome assembly of the insect-repellent terpenoid-producing
601         Lamiaceae species, Callicarpa americana. *GigaScience*, *9*, 1–11.
602         https://doi.org/10.1093/gigascience/giaa093

603     Hamilton, J.P., Vaillancourt, B., Wood, J.C., Wang, H., Jiang, J., Soltis, D.E., Buell, C.R., & Soltis, P.S.
604         (2023). Chromosome-scale genome assembly of the 'Munstead' cultivar of Lavandula
605         angustifolia. *BMC Genomic Data*, *24*, 75

606     Han, M. V., Thomas, G.W.C., Lugo-Martinez, J., & Hahn, M.W. (2013). Estimating gene gain and
607         loss rates in the presence of error in genome assembly and annotation using CAFE 3.
608         *Molecular Biology and Evolution*, *30*, 1987–1997. https://doi.org/10.1093/molbev/mst100

609     Hardigan, M.A., Crisovan, E., Hamilton, J.P., Kim, J., Laimbeer, P., Leisner, C.P., Manrique-
610         Carpintero, N.C., Newton, L., Pham, G.M., Vaillancourt, B., Yang, X., Zeng, Z., Douches, D.S.,
611         Jiang, J., Veilleux, R.E., & Buella, C.R. (2016). Genome reduction uncovers a large dispensable
612         genome and adaptive role for copy number variation in asexually propagated *Solanum*
613         *tuberosum*. *Plant Cell*, *28*, 388–405. https://doi.org/10.1105/tpc.15.00538

He, Z., Feng, X., Chen, Q., Li, L., Li, S., Han, K., Guo, Z., Wang, J., Liu, M., Shi, C., Xu, S., Shao, S., Liu, X., Mao, X., Xie, W., Wang, X., Zhang, R., Li, G., Wu, W., Zheng, Z., Zhong, C., Duke, N.C., Boufford, D.E., Fan, G., Wu, C.I., Ricklefs, R.E., & Shi, S. (2022). Evolution of coastal forests based on a full set of mangrove genomes. *Nature Ecology and Evolution*, *6*, 738–749. https://doi.org/10.1038/s41559-022-01744-9

Healey, A., Furtado, A., Cooper, T., & Henry, R.J. (2014). Protocol: A simple method for extracting next-generation sequencing quality genomic DNA from recalcitrant plant species. *Plant Methods*, *10*, 1–8. https://doi.org/10.1186/1746-4811-10-21/COMMENTS

Hu, G.X., Takano, A., Drew, B.T., Liu, E. De, Soltis, D.E., Soltis, P.S., Peng, H., & Xiang, C.L. (2018). Phylogeny and staminal evolution of *Salvia* (Lamiaceae, Nepetoideae) in East Asia. *Annals of Botany*, *122*, 649–668. https://doi.org/10.1093/AOB/MCY104

Jia, K.H., Liu, H., Zhang, R.G., Xu, J., Zhou, S.S., Jiao, S.Q., Yan, X.M., Tian, X.C., Shi, T. Le, Luo, H., Li, Z.C., Bao, Y.T., Nie, S., Guo, J.F., Porth, I., El-Kassaby, Y.A., Wang, X.R., Chen, C., Van de Peer, Y., Zhao, W., & Mao, J.F. (2021). Chromosome-scale assembly and evolution of the tetraploid Salvia splendens (Lamiaceae) genome. *Horticulture Research*, *8*. https://doi.org/10.1038/s41438-021-00614-y

Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen, J., Mitchell, A., Nuka, G., Pesseat, S., Quinn, A.F., Sangrador-Vegas, A., Scheremetjew, M., Yong, S.-Y., Lopez, R., & Hunter, S. (2014). InterProScan 5: genome-scale protein function classification. *Bioinformatics*, *30*, 1236–1240. https://doi.org/10.1093/bioinformatics/btu031

Kautsar, S.A., Suarez Duran, H.G., Blin, K., Osbourn, A., & Medema, M.H. (2017). plantiSMASH: automated identification, annotation and expression analysis of plant biosynthetic gene clusters. *Nucleic Acids Research*, *45*, W55–W63

Kim, D., Paggi, J.M., Park, C., Bennett, C., & Salzberg, S.L. (2019). Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology*, *37*, 907–915. https://doi.org/10.1038/s41587-019-0201-4

Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H., & Phillippy, A.M. (2017). Canu: Scalable and accurate long-read assembly via adaptive κ-mer weighting and repeat separation. *Genome Research*, *27*, 722–736. https://doi.org/10.1101/gr.215087.116

Kovaka, S., Zimin, A. V., Pertea, G.M., Razaghi, R., Salzberg, S.L., & Pertea, M. (2019). Transcriptome assembly from long-read RNA-seq alignments with StringTie2. *Genome Biology*, *20*, 1–13. https://doi.org/10.1186/s13059-019-1910-1

647    Kumar, S., Suleski, M., Craig, J.M., Kasprowicz, A.E., Sanderford, M., Li, M., Stecher, G., & Hedges,
648        S.B. (2022). TimeTree 5: An Expanded Resource for Species Divergence Times. *Molecular*
649        *Biology and Evolution*, *39*, 1–6. https://doi.org/10.1093/molbev/msac174

650    Lara-Cabrera, S.I., Perez-Garcia, M. de la L., Maya-Lastra, C.A., Montero-Castro, J.C., Godden,
651        G.T., Cibrian-Jaramillo, A., Fisher, A.E., & Porter, J.M. (2021). Phylogenomics of Salvia L.
652        subgenus Calosphace (Lamiaceae). *Frontiers in Plant Science*, *12*.
653        https://doi.org/10.3389/fpls.2021.725900

654    Lee, T.H., Guo, H., Wang, X., Kim, C., & Paterson, A.H. (2014). SNPhylo: A pipeline to construct a
655        phylogenetic tree from huge SNP data. *BMC Genomics*, *15*, 1–6.
656        https://doi.org/10.1186/1471-2164-15-162

657    Li, C.Y., Yang, L., Liu, Y., Xu, Z.G., Gao, J., Huang, Y.B., Xu, J.J., Fan, H., Kong, Y., Wei, Y.K., Hu, W.L.,
658        Wang, L.J., Zhao, Q., Hu, Y.H., Zhang, Y.J., Martin, C., & Chen, X.Y. (2022). The sage genome
659        provides insight into the evolutionary dynamics of diterpene biosynthesis gene cluster in
660        plants. *Cell Reports*, *40*, 111236. https://doi.org/10.1016/j.celrep.2022.111236

661    Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
662        *arXiv preprint*, 1303.3997

663    Li, L., Song, J., Zhang, M., Iqbal, S., Li, Y., Zhang, H., & Zhang, H. (2023). A near complete genome
664        assembly of chia assists in identification of key fatty acid desaturases in developing seeds.
665        *Frontiers in Plant Science*, *14*. https://doi.org/10.3389/fpls.2023.1102715

666    Lichman, B.R., Godden, G.T., Hamilton, J.P., Palmer, L., Kamileen, M.O., Zhao, D., Vaillancourt, B.,
667        Wood, J.C., Sun, M., Kinser, T.J., Henry, L.K., Rodriguez-Lopez, C., Dudareva, N., Soltis, D.E.,
668        Soltis, P.S., Robin Buell, C., & O'Connor, S.E. (2020). The evolutionary origins of the cat
669        attractant nepetalactone in catnip. *Science Advances*, *6*, 1–14.
670        https://doi.org/10.1126/sciadv.aba0721

671    Lovell, J.T., Sreedasyam, A., Schranz, M.E., Wilson, M., Carlson, J.W., Harkess, A., Emms, D.,
672        Goodstein, D.M., & Schmutz, J. (2022). GENESPACE tracks regions of interest and gene copy
673        number variation across multiple genomes. *eLife*, *11*, 1–20.
674        https://doi.org/10.7554/ELIFE.78526

675    Mapleson, D., Accinelli, G.G., Kettleborough, G., Wright, J., & Clavijo, B.J. (2017). KAT: A K-mer
676        analysis toolkit to quality control NGS datasets and genome assemblies. *Bioinformatics*, *33*,
677        574–576. https://doi.org/10.1093/bioinformatics/btw663

678    Marçais, G., Delcher, A.L., Phillippy, A.M., Coston, R., Salzberg, S.L., & Zimin, A. (2018).
679        MUMmer4: A fast and versatile genome alignment system. *PLoS Computational Biology*, *14*,
680        1–14. https://doi.org/10.1371/journal.pcbi.1005944

681    Marçais, G., & Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of
682        occurrences of k-mers. *Bioinformatics*, *27*, 764–770.
683        https://doi.org/10.1093/bioinformatics/btr011

684    Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing
685        reads. *EMBnet.journal*, *17*, 10–12

686    Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A., Sonnhammer, E.L.L., Tosatto,
687        S.C.E., Paladin, L., Raj, S., Richardson, L.J., Finn, R.D., & Bateman, A. (2021). Pfam: The protein
688        families database in 2021. *Nucleic Acids Research*, *49*, D412–D419.
689        https://doi.org/10.1093/nar/gkaa913

690    Mistry, J., Finn, R.D., Eddy, S.R., Bateman, A., & Punta, M. (2013). Challenges in homology search:
691        HMMER3 and convergent evolution of coiled-coil regions. *Nucleic Acids Res*, *41*, e121.
692        https://doi.org/10.1093/nar/gkt263

693    Muñoz, L.A., Cobos, A., Diaz, O., & Aguilera, J.M. (2013). Chia Seed (*Salvia hispanica*): An Ancient
694        Grain and a New Functional Food. *Food Reviews International*, *29*, 394–408.
695        https://doi.org/10.1080/87559129.2013.818014

696    Ou, S., Su, W., Liao, Y., Chougule, K., Agda, J.R.A., Hellinga, A.J., Lugo, C.S.B., Elliott, T.A., Ware,
697        D., Peterson, T., Jiang, N., Hirsch, C.N., & Hufford, M.B. (2019). Benchmarking transposable
698        element annotation methods for creation of a streamlined, comprehensive pipeline.
699        *Genome Biology*, *20*, 1–18. https://doi.org/10.1186/s13059-019-1905-y

700    Pacfici Biosciences. SMRT tools.

701    Pan, X., Chang, Y., Li, C., Qiu, X., Cui, X., Meng, F., Zhang, S., Li, X., & Lu, S. (2023). Chromosome-
702        level genome assembly of Salvia miltiorrhiza with orange roots uncovers the role of
703        Sm2OGD3 in catalyzing 15,16-dehydrogenation of tanshinones. *Horticulture Research*, *10*.
704        https://doi.org/10.1093/hr/uhad069

705    Peláez, P., Orona-Tamayo, D., Montes-Hernández, S., Valverde, M.E., Paredes-López, O., &
706        Cibrián-Jaramillo, A. (2019). Comparative transcriptome analysis of cultivated and wild
707        seeds of *Salvia hispanica* (chia). *Scientific Reports*, *9*, 1–11. https://doi.org/10.1038/s41598-
708        019-45895-5

709    Picard toolkit. (2019). . *Broad Institute, GitHub repository*,

710    Purcell, S., & Chang, C. PLINK 2.0 alpha 2. 3

711    Ranallo-Benavidez, T.R., Jaron, K.S., & Schatz, M.C. (2020). GenomeScope 2.0 and Smudgeplot
712         for reference-free profiling of polyploid genomes. *Nature Communications*, *11*.
713         https://doi.org/10.1038/s41467-020-14998-3

714    Roach, M.J., Schmidt, S.A., & Borneman, A.R. (2018). Purge Haplotigs: allelic contig reassignment
715         for third-gen diploid genome assemblies.. *BMC bioinformatics*, *19*, 460.
716         https://doi.org/10.1186/s12859-018-2485-7

717    Roberts, A., Trapnell, C., Donaghey, J., Rinn, J.L., & Pachter, L. (2011). Improving RNA-Seq
718         expression estimates by correcting for fragment bias. *Genome Biology*, *12*.
719         https://doi.org/10.1186/gb-2011-12-3-r22

720    Shen, Y., Li, W., Zeng, Y., Li, Z., Chen, Y., Zhang, J., Zhao, H., Feng, L., Ma, D., Mo, X., Ouyang, P.,
721         Huang, L., Wang, Z., Jiao, Y., & Wang, H. bin. (2022). Chromosome-level and haplotype-
722         resolved genome provides insight into the tetraploid hybrid origin of patchouli. *Nature
723         Communications*, *13*, 1–15. https://doi.org/10.1038/s41467-022-31121-w

724    Sievers, F., Wilm, A., Dineen, D., Gibson, T.J., Karplus, K., Li, W., Lopez, R., McWilliam, H.,
725         Remmert, M., Söding, J., Thompson, J.D., & Higgins, D.G. (2011). Fast, scalable generation of
726         high-quality protein multiple sequence alignments using Clustal Omega. *Molecular Systems
727         Biology*, *7*. https://doi.org/10.1038/msb.2011.75

728    Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E.M. (2015). BUSCO:
729         assessing genome assembly and annotation completeness with single-copy orthologs.
730         *Bioinformatics*, *31*, 3210–3212. https://doi.org/10.1093/bioinformatics/btv351

731    Smit, A., Hubley, R., & Green, P. *RepeatMasker Open-4.0*.

732    Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
733         phylogenies.              *Bioinformatics*,              *30*,              1312–1313.
734         https://doi.org/10.1093/bioinformatics/btu033

735    Sun, M., Zhang, Y., Zhu, L., Liu, N., Bai, H., Sun, G., Zhang, J., & Shi, L. (2022). Chromosome-level
736         assembly and analysis of the *Thymus* genome provide insights into glandular secretory
737         trichome formation and monoterpenoid biosynthesis in thyme. *Plant Communications*, *3*,
738         100413. https://doi.org/10.1016/j.xplc.2022.100413

739    Tamura, K., Sakamoto, M., Tanizawa, Y., Mochizuki, T., & Bono, H. (2022). Resource Article :
740         Genomes Explored A highly contiguous genome assembly of red perilla (Perilla frutescens)
741         domesticated in Japan 1–8

742  Valdivia-López, M.Á., & Tecante, A. (2015). Chia (*Salvia hispanica*): A Review of Native Mexican
743      Seed and its Nutritional and Functional Properties. *Advances in Food and Nutrition Research*
744      (pp. 53–75). Academic Press Inc.

745  Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C.A., Zeng, Q.,
746      Wortman, J., Young, S.K., & Earl, A.M. (2014). Pilon: An Integrated Tool for Comprehensive
747      Microbial Variant Detection and Genome Assembly Improvement. *PLoS ONE*, *9*, e112963.
748      https://doi.org/10.1371/journal.pone.0112963

749  Wang, L., Lee, M., Sun, F., Song, Z., Yang, Z., & Yue, G.H. (2022). A chromosome-level genome
750      assembly of chia provides insights into high omega-3 content and coat color variation of its
751      seeds. *Plant Communications*, *3*, 100326. https://doi.org/10.1016/j.xplc.2022.100326

752  Wang, Y., Tang, H., Debarry, J.D., Tan, X., Li, J., Wang, X., Lee, T.H., Jin, H., Marler, B., Guo, H.,
753      Kissinger, J.C., & Paterson, A.H. (2012). MCScanX: a toolkit for detection and evolutionary
754      analysis of gene synteny and collinearity. *Nucleic Acids Res*, *40*, e49.
755      https://doi.org/10.1093/nar/gkr1293

756  Xu, Z., Gao, R., Pu, X., Xu, R., Wang, J., Zheng, S., Zeng, Y., Chen, J., He, C., & Song, J. (2020).
757      Comparative Genome Analysis of *Scutellaria baicalensis* and *Scutellaria barbata* Reveals the
758      Evolution of Active Flavonoid Biosynthesis. *Genomics, Proteomics and Bioinformatics*, *18*,
759      230–240. https://doi.org/10.1016/j.gpb.2020.06.002

760  Zhang, Y., Shen, Q., Leng, L., Zhang, D., Chen, S., Shi, Y., Ning, Z., & Chen, S. (2021). Incipient
761      diploidization of the medicinal plant Perilla within 10,000 years. *Nature Communications*,
762      *12*, 1–13. https://doi.org/10.1038/s41467-021-25681-6

763  Zhao, D., Hamilton, J., Bhat, W., Johnson, S., Godden, G., Kinser, T., Boachon, B., Dudareva, D.,
764      Soltis, P., & Soltis, D. (2019)(a). A chromosomal-scale genome assembly of *Tectona grandis*
765      reveals the importance of tandem gene duplication and enables discovery of genes in
766      natural product biosynthetic pathways. *GigaScience*, *8*, giz005

767  Zhao, Q., Yang, J., Cui, M.Y., Liu, J., Fang, Y., Yan, M., Qiu, W., Shang, H., Xu, Z., Yidiresi, R., Weng,
768      J.K., Pluskal, T., Vigouroux, M., Steuernagel, B., Wei, Y., Yang, L., Hu, Y., Chen, X.Y., & Martin,
769      C. (2019)(b). The Reference Genome Sequence of *Scutellaria baicalensis* Provides Insights
770      into the Evolution of Wogonin Biosynthesis. *Molecular Plant*, *12*, 935–950.
771      https://doi.org/10.1016/j.molp.2019.04.002
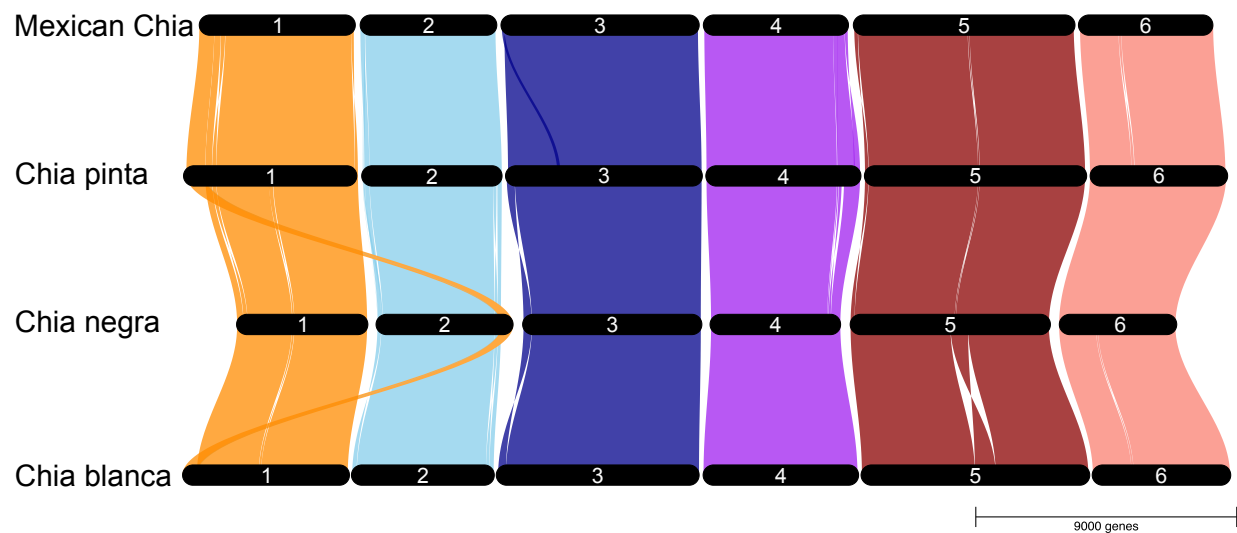
772

773

774    **Figures and Tables**



**Figure 1. Synteny of the Chia genomes.** The top track is the Mexican Chia genome (Alejo-Jacuinde et al., 2023), the second track is the Chia pinta genome reported in this study, the third genome is Chia negra (Wang et al., 2022), and the bottom track represents the Chia blanca genome (Li et al., 2023). The ribbons indicate syntenic blocks between the genomes identified using GENESPACE (v.1.1.10;Lovell et al., 2022).
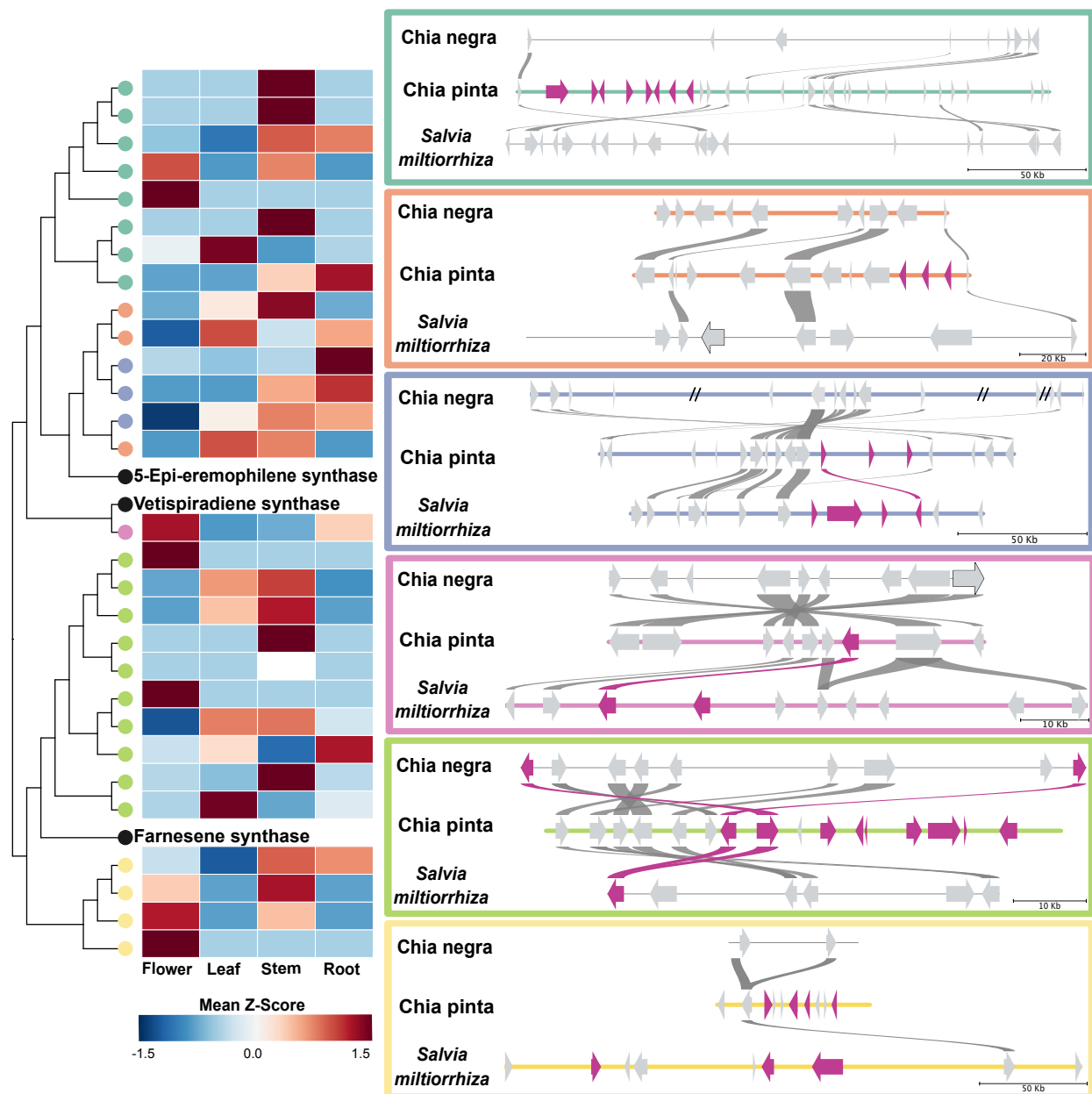
781



**Figure 2. Chia pinta TPS-a Biosynthetic Gene Cluster Expression and Synteny.** A phylogeny of the Chia pinta terpene synthase (TPS-a) genes present in biosynthetic gene clusters (BGCs) with representative functionally characterized reference TPSs is shown. The Chia pinta phylogeny was generated using RAxML (v8.2.12; Stamatakis, 2014). The heatmap of gene expression was constructed from flower, leaf, stem, and root tissue using expression values generated by Cufflinks (v.2.2.1; Roberts et al., 2011) with z-scores range from -1.5 to 1.5. Chia pinta genes (circles on the phylogeny) are colored by their respective BGC and correspond to the outlined syntenic BGCs; genes in black are known TPS. Biosynthetic gene clusters (BGCs) were identified by PlantiSmash (Kautsar et al., 2017) with boxes colored to match the clades in the phylogeny. Syntenic regions were determined using MCScanX (Wang et al., 2012) between Chia pinta, Chia

26

792    negra, and *S. miltorrhiza*. Synteny is indicated as lines between the genes (arrows). The color of

793    the gene and syntenic line is determined by the presumed identity assigned by PlantiSmash

794    where hot pink indicate TPSs; slashes through the line indicate gaps in the assembly. Grey

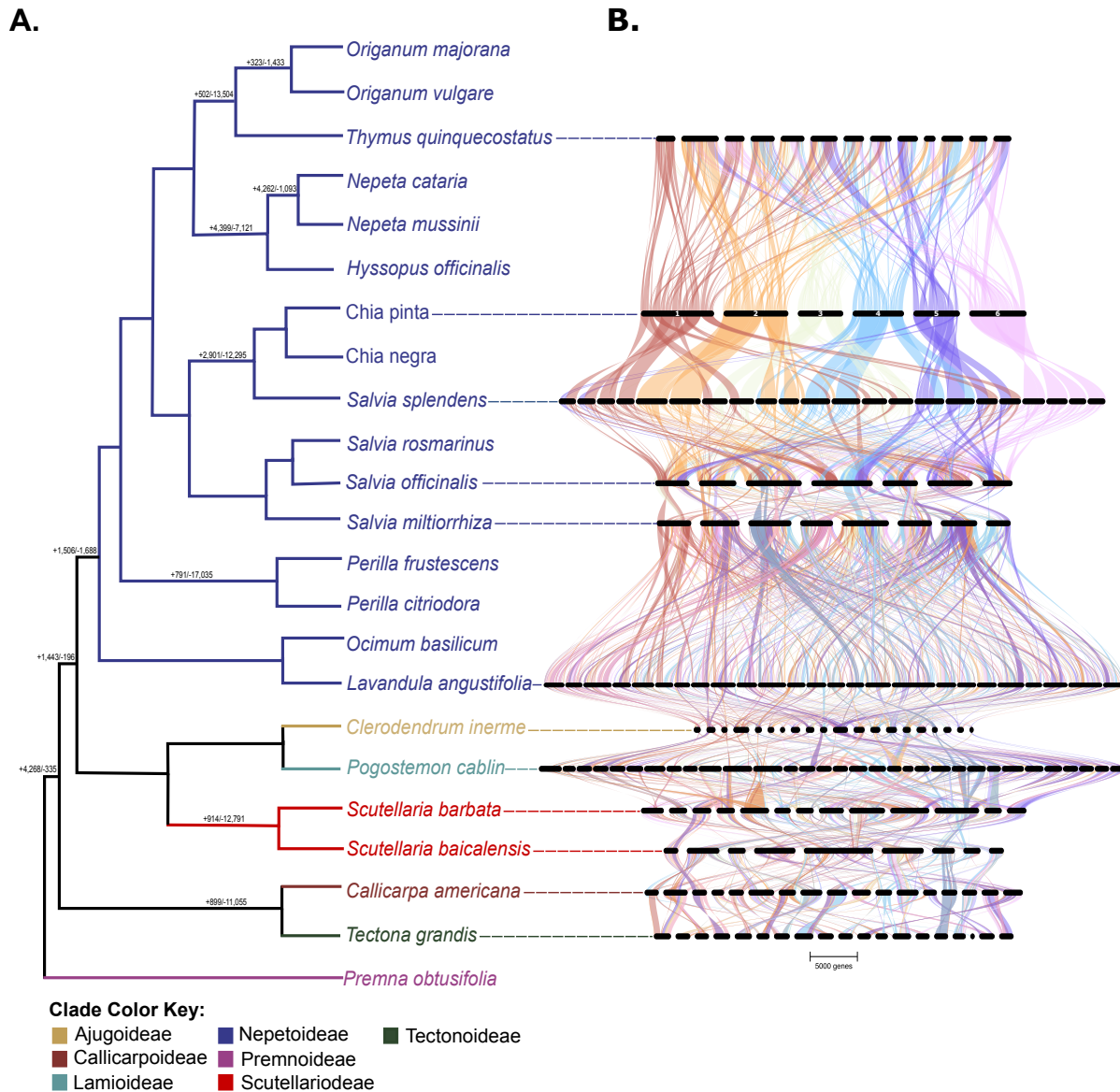795    genome lines indicate that it is not a TPS BGC.

796

**Figure 3. Lamiaceae phylogeny and synteny. A.** A species phylogeny was generated using OrthoFinder (v.2.5.4;Emms & Kelly, 2019) using publicly available chromosome-scale Lamiaceae genomes. Numbers on branches indicated with (+) are gene family expansions and (-) are gene family contractions using CAFE (v.4.2.1; Han et al., 2013). **B.** The GENESPACE (v.1.1.10; Lovell et al., 2022) syntenic map of orthologous regions within chromosome-scale Lamiaceae genome assemblies are shown using the Chia pinta as the reference genome. Chromosomes are scaled by their physical length.
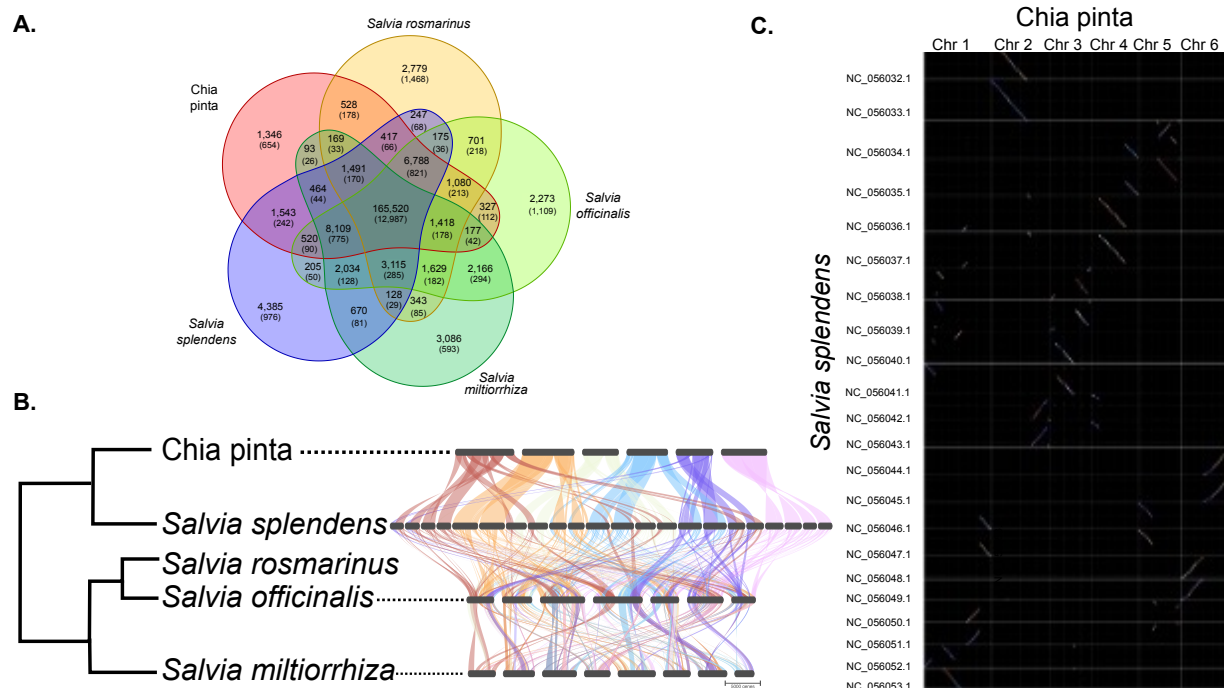
28

**Figure 4. *Salvia* gene orthology and synteny. A.** *Salvia* orthogroup intersections between Chia pinta, *Salvia rosmarinus*, *Salvia officinalis*, *Salvia splendens,* and *Salvia miltiorrhiza* as determined by OrthoFinder (v.2.5.4; Emms & Kelly, 2019). Numbers of orthologous groups and genes in parentheses are reported. **B.** GENESPACE (v.1.1.10; Lovell et al., 2022) syntenic map of orthologous regions within chromosome-scale *Salvia* genome assemblies using Chia pinta as the reference genome. **C.** Synteny dotplot for the anchor genes between Chia pinta and *Salvia splendens* generated in GENESPACE (v.1.1.10; Lovell et al., 2022). Chia pinta includes 21,720 genes with BLAST hits. *Salvia splendens* includes 25,958 genes with blast hits.
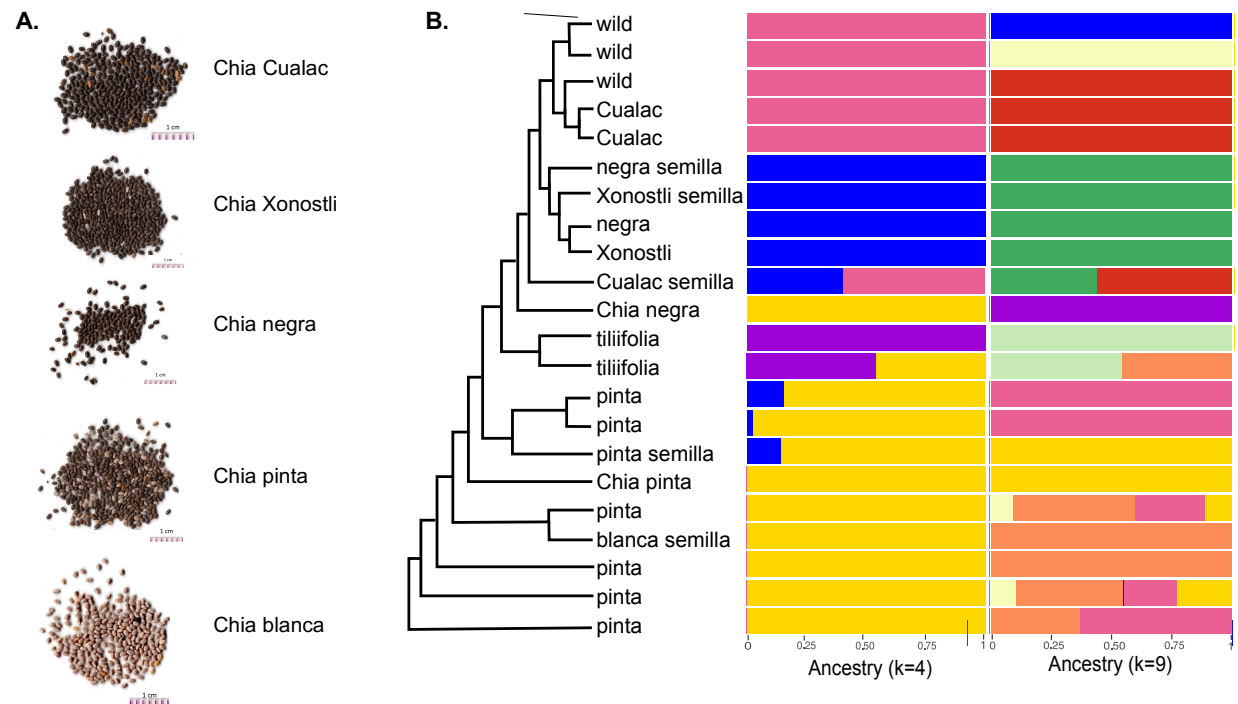
**Figure 5. Population structure of Chia**. **A.** Representative seed images of Chia varieties. **B.** SNP phylogeny was built using SNPhylo (v.20160204; Lee et al., 2014). Admixture (v.1.3.0; Alexander et al., 2009) population structure of 20 Chia accessions and 2 *Salvia tiliifolia* accessions was generated from 156,829 SNPs. Populations from the minima on the cross-validation plot was determined using k=4 and k=9.

**Table 1. Chia pinta Genome Assembly Metrics**

|  | Input assembly | Purged Assembly | Final Chromosome-scale Assembly |
|---|---|---|---|
| **Number of Contigs/ Chromosomes** | 2,094 | 407 | 6 |
| **Total length (bp)** | 425,143,449 | 343,219,856 | 341,980,016 |
| **Maximum Contig Length (bp)** | 9,374,111 | 9,374,111 | 67,233,260 |
| **Minimum Contig Length (bp)** | 1,684 | 2,780 | 57,181,130 |
| **N50 Contig Length (bp)** | 1,150,825 | 1,506,829 | 62,351,092 |
| **Average Contig Length (bp)** | 203,029 | 858,050 | 56,996,669 |

825

826

**Table 2. Chia pinta Genome Annotation Metrics**

| | High Confidence Model Set | High Confidence Representative Model Set | Working Model Set | Working Model Representative Set |
|---|---|---|---|---|
| **Number of Gene Models** | 53,053 | 35,480 | 59,062 | 41,279 |
| **Number of Loci** | 35,480 | 35,480 | 41,279 | 41,279 |
| **Average Transcript Length (bp)** | 3,300.5 | 2,889.0 | 3,104.3 | 2,661.1 |
| **Average CDS Length (bp)** | 1,283.4 | 1,196.6 | 1,216.8 | 1,109.9 |
| **Average Exon Length (bp)** | 280.2 | 283.7 | 279.1 | 280.8 |
| **Average Intron Length (bp)** | 244.2 | 229.8 | 239.8 | 225.2 |
| **Average No. Exons per Model** | 6.1 | 5.3 | 5.8 | 4.9 |
| **Single Exon Transcripts** | 6,105 | 6,043 | 8,062 | 7,999 |

827