

Ten new high-quality genome assemblies for diverse bioenergy sorghum genotypes

1 **William G. Voelker^{1,2*}, Krittika Krishnan^{1,2}, Kapeel Chougule³, Louie C. Alexander Jr.^{1,2},**
2 **Zhenyuan Lu³, Andrew Olson³, Doreen Ware^{3,4}, Kittikun Songsomboon^{1,2}, Cristian**
3 **Ponce^{1,2}, Zachary W. Brenton^{5,6}, J. Lucas Boatwright^{6,7}, Elizabeth A. Cooper^{1,2}**

4 ¹ Dept. of Bioinformatics & Genomics, University of North Carolina at Charlotte, Charlotte, NC
5 USA

6 ²North Carolina Research Campus, Kannapolis, NC USA

7 ³Cold Spring Harbor Research Laboratory, Cold Spring Harbor, NY USA

8 ⁴USDA-ARS NAA, Robert W. Holley Center for Agriculture and Health, Ithaca, NY, USA

9 ⁵Carolina Seed Systems, Darlington, SC USA

10 ⁶Advanced Plant Technology, Clemson University, Clemson, SC USA

11 ⁷Dept. of Plant and Environmental Sciences, Clemson University, Clemson, SC USA

12 *** Correspondence:**

13 William Voelker

14 wvoelker@uncc.edu

15 **Keywords: sorghum, genome assembly, pangenomics, bioenergy, structural variation. (Min.5-**
16 **Max. 8)**

17 **Abstract**

18 Sorghum (*Sorghum bicolor* (L.) Moench) is an agriculturally and economically important staple crop
19 that has immense potential as a bioenergy feedstock due to its relatively high productivity on
20 marginal lands. To capitalize on and further improve sorghum as a potential source of sustainable
21 biofuel, it is essential to understand the genomic mechanisms underlying complex traits related to
22 yield, composition, and environmental adaptations. Expanding on a recently developed mapping
23 population, we generated *de novo* genome assemblies for 10 parental genotypes from this population
24 and identified a comprehensive set of over 24 thousand large structural variants (SVs) and over 10.5
25 million single nucleotide polymorphisms (SNPs). These resources can be integrated into both ongoing
26 and future mapping and trait discovery for sorghum and its myriad uses including food, feed,
27 bioenergy, and increasingly as a carbon dioxide removal mechanism. We show that SVs and
28 nonsynonymous SNPs are enriched in different gene categories, emphasizing the need for long read
29 sequencing in crop species to identify novel variation. Furthermore, we highlight SVs and SNPs
30 occurring in genes and pathways with known associations to critical bioenergy-related phenotypes
31 and characterize the landscape of genetic differences between sweet and cellulosic genotypes.

32 **Introduction**

33 Sorghum (*Sorghum bicolor* (L.) Moench) is a versatile, adaptable, and widely grown cereal crop that
34 is valued for its efficiency, drought tolerance, and ability to grow in marginalized soils (Wayne Smith
35 and Frederiksen, 2000). Present-day genotypes exhibit extensive genetic, phenotypic, morphological,
36 and physiological diversity which stems both from their historical spread and modern breeding

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

37 efforts aimed at optimizing sorghum for different end uses. With its wealth of naturally occurring
38 genetic diversity and advantageous traits, sorghum has enormous value as a sustainable, fast-
39 growing, and high-yielding bioenergy crop (Calviño and Messing, 2012).

40
41 Currently, sorghum is classified into four major ideotypes: grain, sweet, cellulosic, and forage. All of
42 these types can be used in different bioenergy production methods (Wu *et al.*, 2010), but to fully
43 capitalize on their potential, it is essential to gain a better understanding of the genomic changes
44 driving traits related to yield, carbon partitioning, and local adaptation. However, these types of traits
45 are often difficult to dissect due to the nature of their underlying genetic architecture (Brachi, Morris
46 and Borevitz, 2011), which can involve hundreds to thousands of genes and complex mutations that
47 are not easily captured by short-read sequencing.

48
49 Structural genomic mutations are an important source of variation in many species, and can play key
50 roles in phenotypic diversification and evolution. Advances in sequencing technology, especially the
51 advent of high-throughput long-read sequencing, have made the detection of structural variants
52 feasible in many plant species where these types of changes were previously uncharacterized. More
53 recently, there has also been a surge in the generation of pan-genomic data for a number of important
54 crop species, which has offered exciting new insights into the extensive diversity of these plants and
55 the potential influence of complex structural mutations on agronomically important phenotypes
56 (Golicz, Batley and Edwards, 2016; Zhang *et al.*, 2019; Danilevicz *et al.*, 2020; Zhou *et al.*, 2020,
57 2022; Della Coletta *et al.*, 2021; Hufford *et al.*, 2021; Li *et al.*, 2021).

58
59 Previous genomic work in sorghum has linked structural mutations to a number of key traits
60 including dwarfing (Multani *et al.*, 2003), juicy stalks (Zhang *et al.*, 2018), chilling tolerance (Wu *et al.*
61 *et al.*, 2019), and flowering time (Li *et al.*, 2018). A whole-genome comparison of the sweet sorghum
62 genotype ‘Rio’ with ‘BTx623,’ (a short-statured, early maturing grain sorghum) found hundreds of
63 gene presence/absence variations (PAVs), several of which occurred among known sucrose
64 transporters (Cooper *et al.*, 2019). Furthermore, a genome-wide association study (GWAS) exploring
65 the genetic architecture of bioenergy-related traits found that a large deletion in a sorghum-specific
66 iron transporter was linked to stalk sugar accumulation (Brenton *et al.*, 2016, 2020). Most recently,
67 we undertook a broad survey of genome-wide deletions in a panel of nearly 350 diverse sorghum
68 accessions, and found large deletions in multiple genes related to biotic and abiotic stress responses
69 that were unique to particular geographic origins, and appeared to play a role in local adaptation
70 (Songsomboon *et al.*, 2021).

71
72 Taken together, these results suggest that unraveling complex traits in sorghum and other crops will
73 require a comprehensive picture of both structural and single nucleotide mutations. In this study, we
74 have expanded on the recently published Carbon-Partitioning Nested Association Mapping (CP-
75 NAM) population that was developed and publicly released as a key genetic resource for the
76 characterization and improvement of sorghum for multiple different end uses (Boatwright *et al.*,
77 2021, 2022; Kumar *et al.*, 2022). We generated high-quality *de novo* genome assemblies for 10 of the
78 CP-NAM parents and used these genomes to identify millions of novel variants, including a number
79 of large structural variants (SVs) occurring in genes or pathways that could be essential for
80 optimizing sorghum as a bioenergy feedstock.

81 82 **Materials and Methods**

83 *Sample Collection and Sequencing*

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

84 Seeds for each genotype were ordered from the U.S. Department of Agriculture's Germplasm
85 Resource Information Network (GRIN)(<https://www.ars-grin.gov/>) and grown in the greenhouses at
86 the North Carolina Research Campus (NCRC) in Kannapolis, NC. High-molecular-weight DNA was
87 extracted from each sample using a modified high-salt CTAB extraction protocol (Inglis *et al.*, 2018).
88 Purified DNA was sent to the David H. Murdock Research Institute (DHRMI) for quality control,
89 library preparation, and sequencing on a PacBio Sequel I system.

90 *De novo Assembly*

91 Raw subreads for each genotype were combined and converted to FASTQ format using the
92 `bam2fastx` toolkit from PacBio. Reads were then corrected, trimmed, and assembled using
93 Canu(v2.1.1) (Koren *et al.*, 2017). For one of the genotypes, 'Grassl', Canu failed to produce contigs
94 due to reduced read coverage after trimming, so the final assembly was instead produced using
95 Flye(v2.9) with the Canu corrected reads (Kolmogorov *et al.*, 2019).

96 The resulting contigs for all genotypes were scaffolded into chromosomes using RagTag
97 (v2.1.0)(Alonge *et al.*, 2021) and the parameters '-r -g 1 -m 10000000'. Contigs were ordered based
98 on their alignment to the BTx623 v3.1 reference genome (Paterson *et al.*, 2009) with minimap2 (Li,
99 2018). RagTag was run *without* the correction step to avoid unnecessary fragmentation of the contigs
100 and unplaced contigs were discarded. Assembled genome metrics were assessed both before and after
101 scaffolding using QUAST(5.2.0) (Gurevich *et al.*, 2013).

102 *Annotation*

103 Protein and non-coding genes were annotated by building a pan-gene working set using
104 representative pan-gene models selected from a comparative analysis of gene family trees from 18
105 Sorghum genomes (McCormick *et al.* 2018; Deschamps *et al.* 2018; Cooper *et al.* 2019; Wang *et al.*
106 2021; Tao *et al.* 2021) sourced from SorghumBase(<https://www.sorghumbase.org/>). This pan-gene
107 representative was propagated onto the 10 sorghum genome assemblies using Liftoff
108 (v1.6.3)(Shumate and Salzberg 2021) with default parameters. The gene structures were updated with
109 available transcriptome evidence from Btx623 using PASA (v2.4.1)(Haas *et al.* 2003). Additional
110 improvements to structural annotations were done in PASA using full length sequenced cDNAs and
111 sorghum ESTs downloaded from NCBI using the query (EST[Keyword]) AND sorghum[Organism].
112 The working set was assigned Annotation Edit Distance(AED) scores using MAKER-P
113 (v3.0)(Campbell *et al.* 2014) and transcripts with AED score < 1 were classified as protein coding.
114 Those with AED=1 were further filtered to keep any non-BTx623 based models with a minimum
115 protein length of 50 amino acids and a complete CDS as protein coding. The remaining models with
116 AED=1 were classified as non-coding. Gene ID assignment was made as per the existing
117 nomenclature schema established for Sorghum reference genomes(McCormick *et al.* 2018).

118
119 On average, approximately 55 thousand working sets of models were generated for each sorghum
120 line, out of which an average of 41 thousand were coding and roughly 13 thousand were non-coding
121 (Supplementary Table 1). More than half (61%) of the protein coding models mapped to a BTx623
122 reference gene, along with 23% of the non-coding models (Supplementary Figure 1). Functional
123 domain identification was completed with InterProScan (v5.38-76.0) (Jones *et al.* 2014). TRaCE
124 (Olson and Ware 2020) was used to assign canonical transcripts based on domain coverage, protein

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

125 length, and similarity to transcripts assembled by Stringtie. Finally, the protein coding annotations
126 were imported to Ensembl core databases, verified, and validated for translation using the Ensembl
127 API (Stabenau *et al.* 2004).

128

129 In order to assign gene ages, protein sequences were aligned to the canonical translations of gene
130 models from *Zea mays*, *Oryza sativa*, *Brachypodium distachyon*, and *Arabidopsis thaliana* obtained
131 from Gramene release 62 (Tello-Ruiz *et al.* 2020) using USEARCH v11.0.667_i86linux32 (Edgar
132 2010). If there was a hit with minimum sequence identity of 50% (-id 0.5) to an *Arabidopsis* protein,
133 the gene was classified as being from Viridiplanteae, if there was a hit to rice the gene was classified
134 as Poaceae, and if a hit was to maize the gene was classified as Andropogoneae. If there were no hits
135 then the gene was classified as sorghum specific.

136 *Repeat Analysis*

137 Transposable elements (TEs) were identified and annotated in each genome using EDTA (Ou *et al.*,
138 2019). TE-greedy-nester (Lexa *et al.*, 2020) was used to further annotate both complete and
139 fragmented Long Terminal Repeat (LTR) retrotransposons. Sequence divergence in the LTR regions
140 was used to estimate retrotransposon age (SanMiguel *et al.*, 1998; Jedlicka, Lexa and Kejnovsky,
141 2020). The left and right LTR sequences were extracted from the assembled genomes using the
142 coordinates reported by TE-greedy-nester and the `getfasta` tool from the BEDTools
143 package(v2.29.0) (Quinlan and Hall, 2010). For each TE, the two LTR sequences were aligned using
144 Clustal-W (Thompson, Higgins and Gibson, 1994) as implemented in the R package `msa`
145 (Bodenhofer *et al.*, 2015). Genetic distance was calculated based on the K80 model using the
146 `dist.dna` function in the R package `phangorn` (Schliep, 2011). The time of divergence was
147 calculated based on the equation $T=K/(2 * r)$ (Bowen and McDonald, 2001), where T is the time of
148 divergence, K is the genetic distance, and r is the substitution rate. A value of 0.013 mutations per
149 million years was used for r, consistent with the molecular clock rate for LTRs estimated in rice (Ma
150 and Bennetzen, 2004).

151 *Variant Calling*

152 Filtered and scaffolded reads were realigned to the BTx623 reference genome using the `nucmer`
153 program from the MUMmer(v4.0) package (Delcher, Salzberg and Phillippy, 2003; Marçais *et al.*,
154 2018) with the following parameters '-c 100 -b 500 -l 50'. Alignments were filtered using the
155 `delta-filter` program from the MUMmer package with the parameters '-m -i 90 -l 100' and
156 converted to coordinate files using `show-coords` with the parameters '-THrd'. Variants were then
157 called using Syri(v1.6)(Goel *et al.*, 2019).

158 Individual Syri VCF files were split by variant type (SNPs, Deletions, Insertions, Inversions, and
159 Translocations) resulting in separate files for each variant type for each genotype. Insertions or
160 deletions smaller than 50 bp were classified as small indels while those equal to or larger than 50 bp
161 were classified as SVs. More complex SV types that could not be validated with raw reads were not
162 considered for further analysis.

163 The Syri program produces a nonstandard VCF format which includes information on variants from
164 overlapping syntenic blocks. This can result in duplicated variants and fragmented insertions that
165 must be addressed before subsequent analysis with downstream tools. Duplicates of existing variants
166 were removed for all variant types, and fragmented insertions were combined into single variants

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

167 (Supplementary Figure 2). These processed variant files were then zipped and indexed using `bgzip`
168 and `tabix` (Li *et al.*, 2009) and then merged across genotypes using the merge function from the
169 `bcftools` package with the parameters ‘-O -I ‘ChrB:join,Parent:join,DupType:join,modified:join’ -O
170 v’. This resulted in one variant file for each type of variant that included the genotypes for all
171 individuals. Insertions, deletions, and SNPs were then annotated using SIFT (v2.4)(Vaser *et al.*,
172 2016) and the BTx623 version 3.1.1 annotation to identify overlap with genes for insertions and
173 deletions and missense prediction for single nucleotide variants.

174 *Phylogeny*

175 Gene PAVs were called from pan-gene lift-off annotation information using custom python scripts.
176 PAVs for each genotype were encoded as a binary vector (with 0 indicating gene absence, and 1
177 indicating presence). Distance between genotypes was then calculated using the `dist()` function from
178 the `stats(v3.6.2)` package in R using the Jaccard distance, and a phylogenetic tree was constructed
179 using the `NJ()` function from the `phangorn` package.

180 *Gene Ontology Analysis*

181 Gene ontology (GO) terms for genes affected by large insertions and deletions or nonsynonymous
182 SNPs were curated from the publicly available annotation information file associated with BTx623
183 v3.1.1 in phytozome (<https://phytozome-next.jgi.doe.gov/>). GO enrichment analysis was performed
184 using the R package `topGO(v1.0)` (Alexa and Rahnenfuhrer, 2016). The classic Fisher’s Test was
185 used to assess significance of enriched terms, and terms with a p-value <0.05 were considered
186 significant and kept for further analysis. Redundant and highly similar GO terms were defined and
187 reduced based on semantic similarity using the R packages `AnnotationForge` (Carlson and Pages,
188 2022) and `rrvgo` (Sayols, 2020).

189 **Results**

190 *Assembly Quality and Characteristics*

191 To capture the genetic diversity of bioenergy sorghum, we sequenced the parents of the previously
192 established CP-NAM population, which included globally diverse genotypes representative of sweet,
193 cellulosic, grain and forage type bioenergy sorghums (Boatwright *et al.*, 2021)(Table 1). The initial
194 contig-level assemblies showed a range of N50 values, with the lowest being 176 kb and the highest
195 at over 3 Mbp (Supplementary Table 2). The three sweet genotypes in particular had a higher number
196 of raw reads and more contiguous assemblies than the other types (Figures 1A and 1B), most likely
197 as a result of differences in the effectiveness of the extraction protocol. After scaffolding and filtering
198 unplaced contigs, all 10 genotypes showed similar levels of high contiguity, with final assembly sizes
199 that were 90-98% the size of the BTx623 reference genome and over 90% of known BTx623 genes
200 contained within the scaffolds (Figures 1C and 1D).

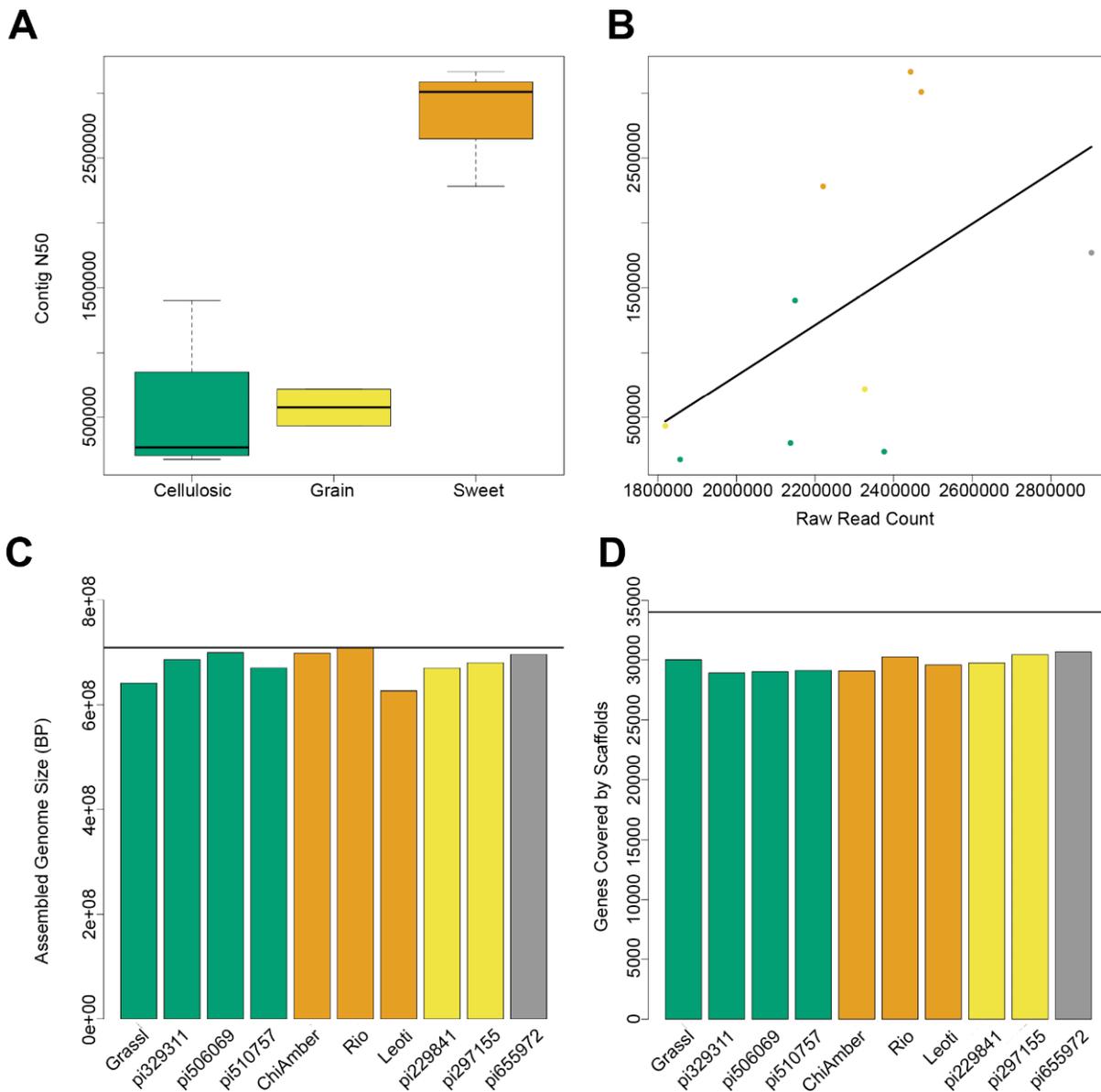
201

202

203

204

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes



205

206 Figure 1. Assembly metrics for 10 sorghum genotypes. A) Contig N50 levels for different ideotypes show higher
207 contiguity for sweet genotypes. B) Raw read counts prior to assembly are highly correlated with contig N50, and sweet
208 genotypes (orange) have higher read counts than cellulosic (green) or grain (yellow) genotypes. C) Assembled genome
209 size after scaffolding and filtering for each genotype shows that despite differences in mean contig size, the final
210 assemblies for both sweet and non-sweet types are very close to the expected reference genome size (horizontal black
211 line). D) The number of BTx623 genes contained within the final scaffolds is very similar across all genotypes regardless
212 of type.

213

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

214 Table 1. Genotype Origins, Races, and Types.

Name	Alternate ID	Race	Origin	Type
Grassl	PI 154844	Caudatum	Uganda	Sweet & Cellulosic
PI 329311	IS 11069	Durra	Ethiopia	Cellulosic
PI 506069	Mbonou	Guinea-bicolor	Togo	Cellulosic
PI 510757	AP79-714	Durra	Cameroon	Cellulosic
Chinese Amber	PI 22913	Bicolor	China	Sweet
Rio	PI 563295	Durra-caudatum	USA	Sweet
Leoti	PI 586454	Kafir-bicolor	Hungary	Sweet
PI 229841	IS 2382	Kafir	South Africa	Grain
PI 297155	IS 13633	Kafir	Uganda	Grain, Forage
PI 655972	Pink Kafir	Kafir	USA	Forage

215 Information adapted from GRIN and (Boatwright *et al.*, 2021).

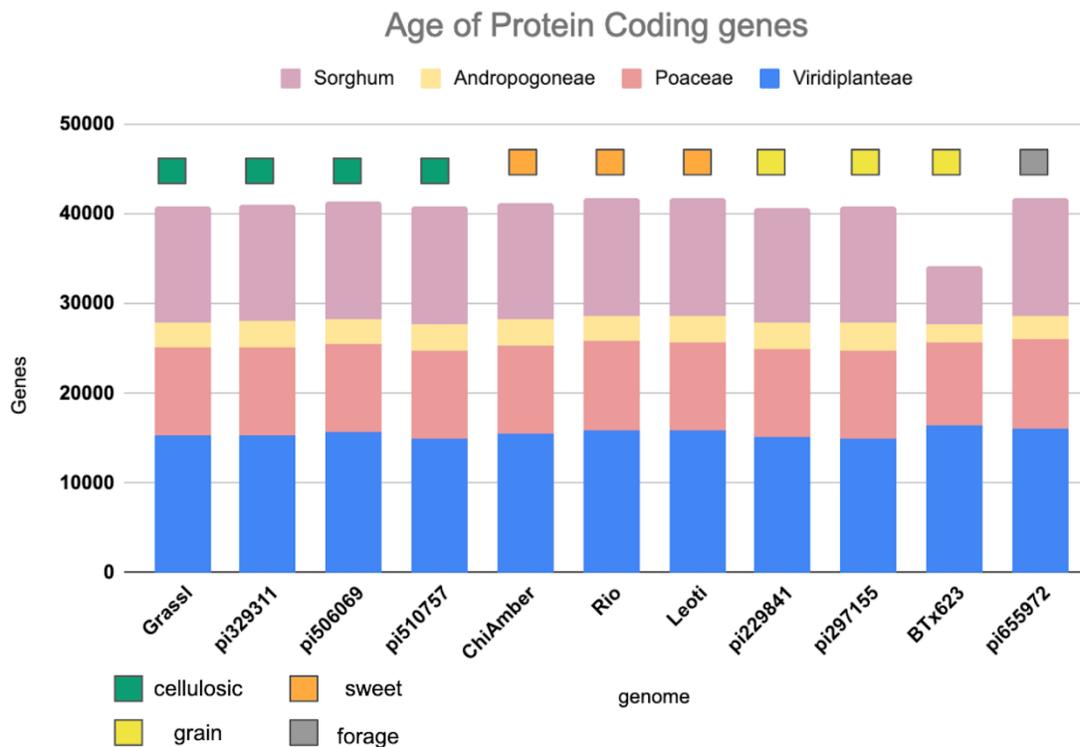


Figure 2. Age of protein coding genes among the sorghum lines based on minimum sequence identity. Bar color indicates the level of phylogenetic conservation, with blue indicating genes conserved across monocots and dicots; peach indicating the proportion of genes shared among the grasses; yellow indicating the proportion of genes shared between sorghum and maize, and light purple representing the proportion of sorghum-specific genes.

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

218 *Gene Annotation*

219 Genes shared across deeper evolutionary time scales were more conserved than sorghum-specific
220 genes (Figure 2). The sweet genotypes show slightly more conserved genes when compared to other
221 genotypes (Figure 2). Around 36.69 percent of genes were found to be core to all genotypes, 50.32
222 percent were shell genes (present in more than one genome, but not all of the genomes), and 12.99
223 percent were found to be cloud genes (unique to a single genome) (Supplementary Figure 3). Of shell
224 genes identified, 44 and 45 were identified to be exclusive to all sweet and all non-sweet genotypes
225 respectively.

226 *Genomic Landscape of Variation*

227 Over 10.5 million single nucleotide variants were called across the 10 genomes, as well as over 7.4
228 million small indels and over 24 thousand large structural variants (insertions and deletions ≥ 50 bp)
229 (Figure 3, Tables 2 and 3). Well over half (~65%) of these variants were defined as cloud variation
230 (Table 3), while the remaining variants were mostly shell. Only a small handful of core variants were
231 present in all of the genotypes except the BTx623 reference. Phylogenetic relationships were inferred
232 using gene presence/absence to estimate genetic distance (Supplementary Figure 4), demonstrating
233 that sweet, cellulosic, and grain genotypes come from separate clades within the category of
234 bioenergy-type sorghum.

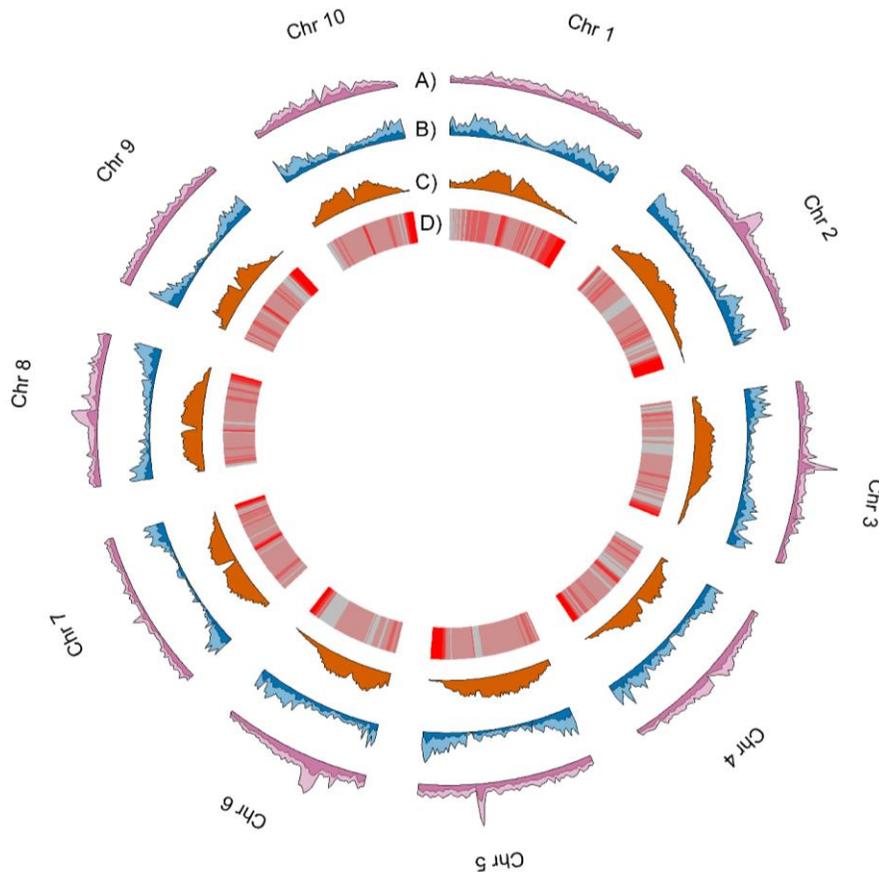


Figure 3. Genomic landscape of variation averaged across the 10 genomes. Density estimates in tracks A-C were performed in 1Mb non-overlapping sliding windows. A) and B) respectively show average SNP density and average SV density, with lighter colors indicating cloud variants and darker colors indicating shell and core variants. C) shows the average TE density, and D) shows TE age averaged across 1Mb sliding windows. Red indicates younger TEs while gray indicates older.

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

235 Table 2. Variants found in each NAM parent genotype.

Genotype	Deletions(bp>=50)	Insertions(bp>=50)	Indels(bp<50)	SNPs	Nonsynonymous
Grassl	2,721	1,714	976,703	2,659,850	37,265
PI 329311	3,560	1,956	1,319,281	3,321,035	47,482
PI 506069	3,531	1,865	888,425	3,003,469	47,555
PI 510757	2,952	1,919	1,593,228	2,859,852	44,168
Chinese Amber	3,560	1,744	994,023	2,975,137	48,780
Rio	2,563	1,791	717,304	2,119,637	35,714
Leoti	3,279	1,435	785,360	2,790,452	43,473
PI 229841	2,830	1,490	1,447,030	2,546,090	41,679
PI 297155	2,412	1,335	1,151,594	2,052,203	34,863
PI 655972	2,401	1,113	631,705	1,953,106	32,758

236

237 Table 3. Core vs. Shell vs. Cloud variants

Type	Deletions	Insertions	Total SVs	Indels	SNPs
Core	34	28	62	12,231	103,065
Shell	6,306	2,250	8,556	1,246,552	5,245,181
Cloud	7,855	8,232	16,087	6,195,713	5,416,344
Total	14,195	10,510	24,705	7,454,496	10,764,590

238

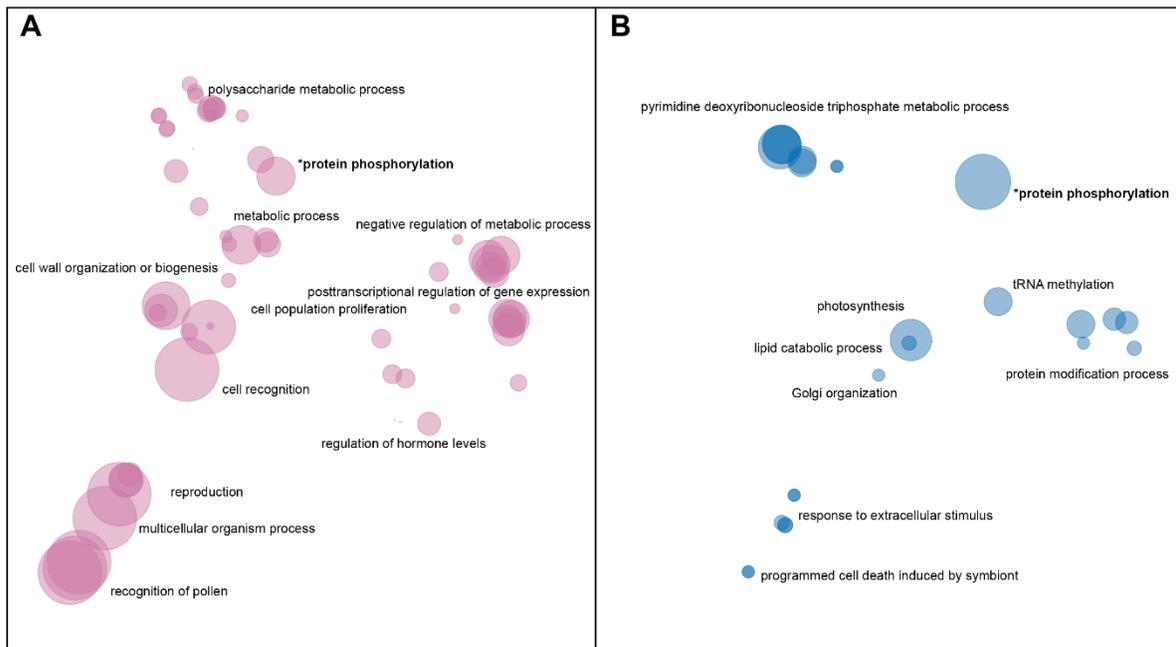
239 *Genes Affected by Structural Variants and SNPs*

240 There were a total of 171,000 SNPs that were found to be both located in genic regions and encoding
 241 nonsynonymous variants, and more than 2.5 thousand large SVs present in genic regions. GO
 242 enrichment analyses of affected genes revealed that SNPs and SVs tended to impact distinct
 243 categories of genes (Figure 4), with protein phosphorylation being the only significant category to
 244 appear in both datasets.

245 In addition to protein phosphorylation, genes impacted by large insertions or deletions showed
 246 enrichment in GO categories related to Golgi vesicle transport, photosynthesis, nucleoside
 247 metabolism, protein modifications, and programmed cell death (Figure 4B). Nonsynonymous SNPs,
 248 on the other hand, were enriched in genes involved in pollen-pistil interactions, cell wall biogenesis,
 249 cell proliferation, posttranscriptional regulation and polysaccharide metabolism (Figure 4A).

250

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes



251

252 Figure 4. Enriched GO terms for genes impacted by A) nonsynonymous SNPs and B) large SVs. GO terms in each
253 dataset were clustered and plotted based on semantic similarity as described in the Materials and Methods. Circle size is
254 proportional to p-value, with larger circles indicating more significant terms.

255 Repeat Analysis

256 Overall the TE composition was highly similar across all 10 genotypes (Figure 5 and 3), with the
257 LTR-Gypsy superfamily comprising the majority of elements. The age analysis revealed an
258 abundance of younger TEs, with a mean age of 1.28 million years old along with a high frequency of
259 very young TEs approximately 0.1 million years old and very few old TEs (6-8 million years)
260 (Figure 5; Supplementary Figure 5). Most (97.5%) of the TEs were non-nested, with TE-greedy-
261 nester reporting the presence of only a handful (2.5%) of nested TEs. The overall distribution of TE
262 age followed a similar pattern across all of the genotypes, with younger TEs being randomly
263 distributed throughout the genome (Figure 3, Supplementary Figure 6A-J) as previously observed by
264 (Paterson *et al.*, 2009).

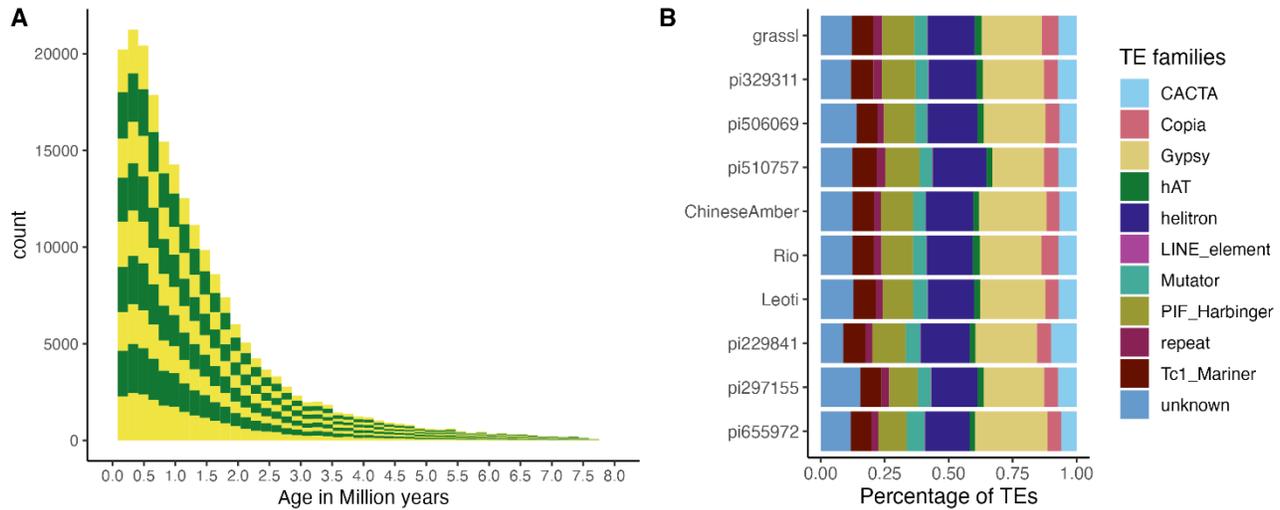
265 Differences in Sweet and Non-Sweet Genotypes

266 Structural variants that were present in all three sweet genotypes (Leoti, ChineseAmber, and Rio) but
267 either absent from or rare among non-sweet genotypes, were significantly enriched among genes with
268 functions related to metal ion transport, in particular iron ion transport, as well as genes involved in
269 oxidative stress response, cell cycle arrest, and phosphatidylserine biosynthetic processes.
270 Conversely, variants found only in all of the non-sweet genotypes tended to impact very different
271 categories of genes, such as those involved in glycolytic processes, cytochrome assembly, and both
272 RNA and DNA regulation (Figure 6).

273

274

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes



275
276

277 Figure 5. TE age and composition. A) The distribution of TE counts by age across all genotypes. Alternating colors
278 indicate different genotypes, with the Y-axis labeling the number of TEs and the X-axis labeling their age in millions of
279 years. B) The proportion of superfamilies of TEs based on average counts of each superfamily across all genomes.

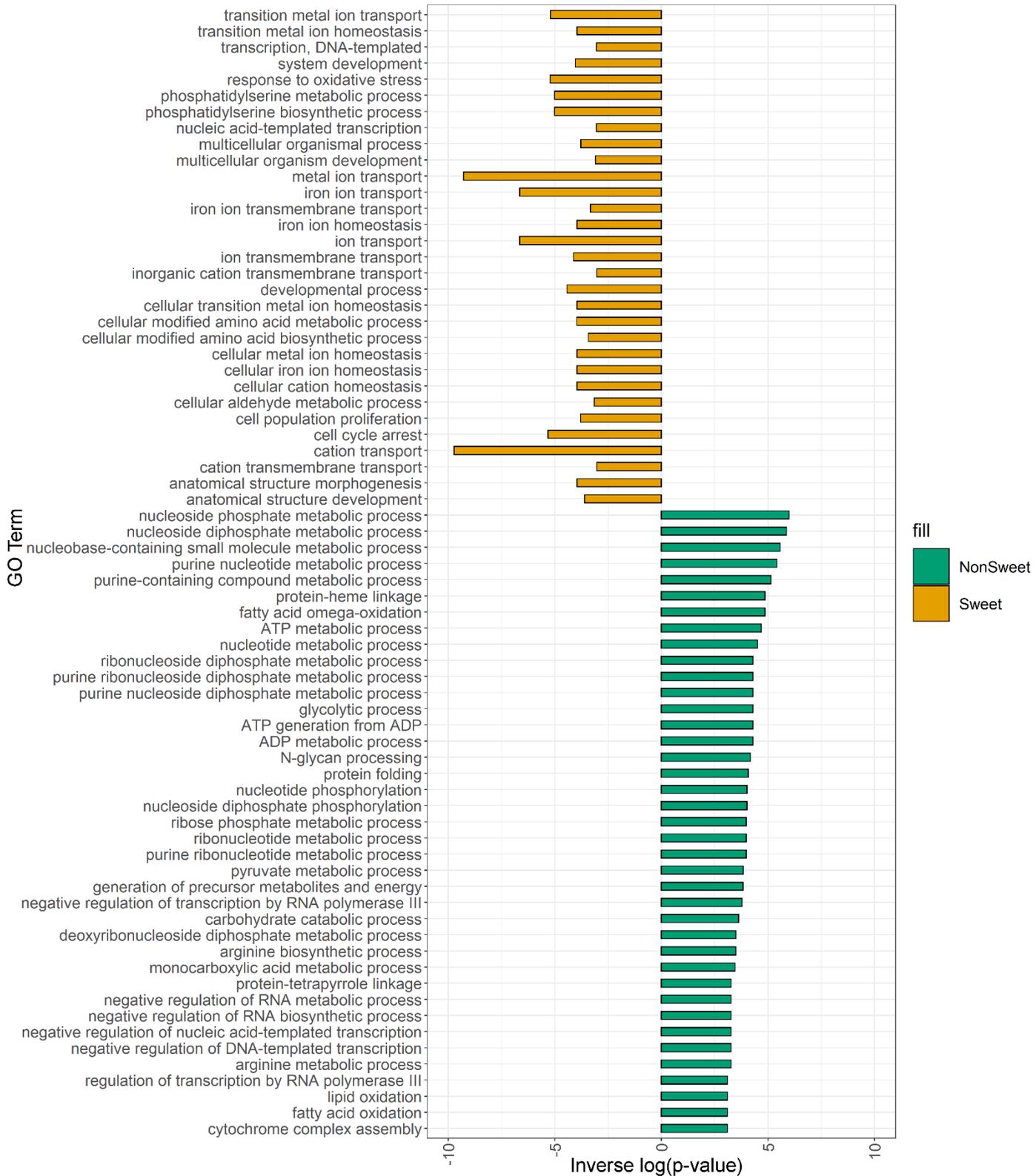
280 Discussion

281 Unraveling the molecular mechanisms controlling complex traits such as carbon partitioning, yield,
282 and stress response is an essential step for crop improvement efforts aimed at creating effective and
283 sustainable bioenergy feedstocks for the future. However, not only do these types of traits often
284 involve changes in large numbers of genes, but an ever-increasing number of pan-genomics studies
285 in crop plants have demonstrated that these changes can encompass complex structural mutations in
286 addition to SNPs (Cooper *et al.*, 2019; Zhang *et al.*, 2019; Brenton *et al.*, 2020; Zhou *et al.*, 2020,
287 2022; Hufford *et al.*, 2021; Songsomboon *et al.*, 2021). Therefore, the development of multiple
288 reference-quality genomes within crop species is critical to the exploration of complex genetic
289 architectures and has clear benefits when compared to a single reference genome, especially in the
290 case of larger structural variants (Della Coletta *et al.*, 2021). By *de novo* assembling 10 new high-
291 quality genomes for the parents of the CP-NAM population (Boatwright *et al.*, 2022), we have been
292 able to uncover millions of novel variants, including thousands of large insertions and deletions.

293 Importantly, we found that SVs within coding regions impacted different types of genes compared to
294 SNPs, highlighting the importance of incorporating both into future trait mapping studies. Many
295 nonsynonymous SNPs that were segregating among the genotypes occurred in gene categories that
296 have previously been linked to carbon allocation in sorghum and other closely related species. For
297 instance, protein phosphorylation induces key signaling cascades in plants that control a variety of
298 processes, and protein kinases have been shown to be highly differentially expressed in both sweet
299 sorghum (Cooper *et al.*, 2019) and sugarcane (Waclawovsky *et al.*, 2010) during stem sugar
300 accumulation. Similarly, genes involved in the regulation of plant hormones such as auxin were also
301 enriched for non-coding SNPs, and these pathways are known to be essential for vegetative plant
302 growth and stem elongation, both of which are key phenotypes for biomass accumulation (Kebrom,
303 McKinley and Mullet, 2017).

304

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes



305

Figure 6. Enriched GO terms for genes impacted by SVs and Indels in both Non-Sweet and Sweet Genotypes. Orange bars indicated gene categories in Sweet genotypes that were significantly impacted ($p < 0.05$). Green bars indicated gene categories in Non-Sweet genotypes that were significantly impacted ($p < 0.05$). The length of each bar corresponds to significance ($-\log(p\text{-value})$).

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

306 Like SNPs, gene-impacting SVs were also found to affect many genes related to protein
307 phosphorylation; in fact, this was the top category among genes containing large variants. But other
308 categories enriched for high-impact insertions and deletions were distinct from the SNP dataset, and
309 contained many genes involved in pathways related to both abiotic and biotic stress responses, which
310 has been observed before in diverse bioenergy sorghums (Songsomboon *et al.*, 2021). Additionally
311 our study identified structural variants affecting genes involved in tRNA nucleoside modifications,
312 programmed cell death in response to symbionts, and photosynthetic light response, all of which
313 were previously identified by other studies as GO terms of interest in relation to sorghum stress
314 response (Ortiz, Hu and Salas Fernandez, 2017; Wang *et al.*, 2017).

315 SVs strictly occurring in either sweet or non-sweet genotypes also offer unique insights into the
316 differences between these types that could be key to dissecting differences in carbon allocation in
317 sorghum. Of particular interest is the fact that SVs restricted to sweet sorghum genotypes affected
318 many genes related to metal metabolism and iron transport. This connection between iron transport
319 and sugar accumulation has been observed in other comparative genomic studies of sorghum
320 (Brenton *et al.*, 2016, 2020; Cooper *et al.*, 2019), and appears to be a key factor distinguishing sweet
321 sorghums from both cellulosic and grain types.

322 Over a third of protein coding genes and over 75 percent of noncoding genes annotated in this study
323 did not map back to the Btx623 reference genome. With a growing number of studies illustrating the
324 importance of noncoding DNA and RNA as potential regulatory elements (Waititu *et al.* 2020), it is
325 evident that large pan-genome annotations are vital in quickly identifying and annotating potential
326 regulatory ‘pseudo-genes’ as well as protein coding genes that are divergent from the common
327 reference. Previous pan-genome studies in sorghum and maize have identified high levels of gene
328 content variation, with 53-64 percent of genes identified as non-core (Tao *et al.*, 2021; Ruperao *et al.*,
329 2021; Hufford *et al.*, 2021). We corroborate these findings with about 63 percent of our genes being
330 identified as either shell or cloud to our population, despite this particular population lacking wild
331 representation, indicating relatively high amounts of latent variation, even among domesticated
332 varieties of sorghum.

333 Taken together, our results demonstrate the value of exploring genome-wide patterns of both SNPs
334 and larger structural variants to gain new insights into the genetic architectures of complex and
335 agronomically important traits. To advance both sorghum breeding efforts and our understanding of
336 crop plant evolution, we have generated this new extensive dataset that is publicly available through
337 SorghumBase (Gladman *et al.*, 2022) and which can be readily integrated into an already valuable
338 genetic resource for future mapping studies.

339 Nomenclature

340 CP-NAM: Carbon Partitioning Nested Association Mapping

341 SV: Structural Variant

342 SNP: Single Nucleotide Polymorphism

343 TE: Transposable Element

344 LTR: Long Terminal Repeat

345 GO: Gene Ontology

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

346 **Conflict of Interest**

347 *The authors declare that the research was conducted in the absence of any commercial or financial*
348 *relationships that could be construed as a potential conflict of interest.*

349 **Author Contributions**

350 WGV: Writing, variant analysis, created figures and tables, performed scaffolding.
351 KK: Performed TE analysis, Alignments, and Variant calling. Wrote corresponding methods
352 sections.
353 LCA: Aided in scripting of figure creation and filtering of variants.
354 KS: Growing and DNA Extraction of plant material.
355 CP: Aided in genome assembly.
356 KC, ZL, AO: Gene and transposable element annotations.
357 DW: Experimental design, writing.
358 ZWB: designed CP-NAM population, provided genetic materials
359 JLB: development and release of CP-NAMs
360 EAC: Writing, created figures, conceived the project, advised, and helped direct analysis.

361 **Funding**

362 This research was supported by startup funds from UNC Charlotte and the United States Department
363 of Agriculture grant USDA-ARS 8062-21000-041-00D.

364 **Acknowledgments**

365 The authors would like to thank S. Kresovich and M. Myers for providing plant materials, J. Lotito
366 and N.C. State for providing and overseeing the greenhouse and growth chambers facilities at the
367 N.C. Research Campus, and the DHMRI Genomics Core for providing sequencing services. The
368 authors would also like to acknowledge the University Research Computing team at UNC Charlotte
369 and S. Blanchard for providing essential IT support and resources.

370 **References**

- 371 Alexa and Rahnenfuhrer (2016) ‘topGO: Enrichment analysis for Gene Ontology. R package version 2.28. 0’,
372 *Cranio: the journal of craniomandibular practice* [Preprint].
- 373 Alonge, M. *et al.* (2021) ‘Automated assembly scaffolding elevates a new tomato system for high-throughput
374 genome editing’, *bioRxiv*. doi:10.1101/2021.11.18.469135.
- 375 Boatwright, J.L. *et al.* (2021) ‘Genetic characterization of a Sorghum bicolor multiparent mapping population
376 emphasizing carbon-partitioning dynamics’, *G3*, 11(4). doi:10.1093/g3journal/jkab060.
- 377 Boatwright, J.L. *et al.* (2022) ‘Dissecting the Genetic Architecture of Carbon Partitioning in Sorghum Using
378 Multiscale Phenotypes’, *Frontiers in plant science*, 13, p. 790005.
- 379 Bodenhofer, U. *et al.* (2015) ‘msa: an R package for multiple sequence alignment’, *Bioinformatics*, 31(24),
380 pp. 3997–3999.
- 381 Bowen, N.J. and McDonald, J.F. (2001) ‘Drosophila euchromatic LTR retrotransposons are much younger
382 than the host species in which they reside’, *Genome research*, 11(9), pp. 1527–1540.

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

- 383 Brachi, B., Morris, G.P. and Borevitz, J.O. (2011) ‘Genome-wide association studies in plants: the missing
384 heritability is in the field’, *Genome biology*, 12(10), p. 232.
- 385 Brenton, Z.W. *et al.* (2016) ‘A Genomic Resource for the Development, Improvement, and Exploitation of
386 Sorghum for Bioenergy’, *Genetics*, 204(1), pp. 21–33.
- 387 Brenton, Z.W. *et al.* (2020) ‘Species-Specific Duplication Event Associated with Elevated Levels of
388 Nonstructural Carbohydrates in Sorghum bicolor’, *G3*, 10(5), pp. 1511–1520.
- 389 Calviño, M. and Messing, J. (2012) ‘Sweet sorghum as a model system for bioenergy crops’, *Current opinion
390 in biotechnology*, 23(3), pp. 323–329.
- 391 Campbell, Michael S., Carson Holt, Barry Moore, and Mark Yandell. 2014. “Genome Annotation and
392 Curation Using MAKER and MAKER-P.” *Current Protocols in Bioinformatics / Editorial Board, Andreas D.
393 Baxevanis ... [et Al.]* 48 (December): 4.11.1–39.
- 394 Carlson and Pages (2022) ‘AnnotationForge: code for building annotation database packages’, *R package
395 version* [Preprint].
- 396 Cooper, E.A. *et al.* (2019) ‘A new reference genome for Sorghum bicolor reveals high levels of sequence
397 similarity between sweet and grain genotypes: implications for the genetics of sugar metabolism’, *BMC
398 genomics*, 20(1), p. 420.
- 399 Danilevicz, M.F. *et al.* (2020) ‘Plant pangenomics: approaches, applications and advancements’, *Current
400 opinion in plant biology*, 54, pp. 18–25.
- 401 Delcher, A.L., Salzberg, S.L. and Phillippy, A.M. (2003) ‘Using MUMmer to identify similar regions in large
402 sequence sets’, *Current protocols in bioinformatics / editorial board, Andreas D. Baxevanis ... [et al.]*, Chapter
403 10, p. Unit 10.3.
- 404 Della Coletta, R. *et al.* (2021) ‘How the pan-genome is changing crop genomics and improvement’, *Genome
405 biology*, 22(1), p. 3.
- 406 Deschamps, Stéphane, Yun Zhang, Victor Llaca, Liang Ye, Abhijit Sanyal, Matthew King, Gregory May, and
407 Haining Lin. 2018. “A Chromosome-Scale Assembly of the Sorghum Genome Using Nanopore Sequencing
408 and Optical Mapping.” *Nature Communications* 9 (1): 4844.
- 409 Gladman, N. *et al.* (2022) ‘SorghumBase: a web-based portal for sorghum genetic information and community
410 advancement’, *Planta*, 255(2), p. 35.
- 411 Goel, M. *et al.* (2019) ‘SyRI: finding genomic rearrangements and local sequence differences from whole-
412 genome assemblies’, *Genome biology*, 20(1), p. 277.
- 413 Golicz, A.A., Batley, J. and Edwards, D. (2016) ‘Towards plant pangenomics’, *Plant biotechnology journal*,
414 14(4), pp. 1099–1105.
- 415 Gurevich, A. *et al.* (2013) ‘QUAST: quality assessment tool for genome assemblies’, *Bioinformatics*, 29(8),
416 pp. 1072–1075.
- 417 Haas, Brian J., Arthur L. Delcher, Stephen M. Mount, Jennifer R. Wortman, Roger K. Smith Jr, Linda I.
418 Hannick, Rama Maiti, *et al.* (2003). “Improving the Arabidopsis Genome Annotation Using Maximal
419 Transcript Alignment Assemblies.” *Nucleic Acids Research* 31 (19): 5654–66.
- 420 Hufford, M.B. *et al.* (2021) ‘De novo assembly, annotation, and comparative analysis of 26 diverse maize
421 genomes’, *Science*, 373(6555), pp. 655–662.

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

- 422 Inglis, P.W. *et al.* (2018) ‘Fast and inexpensive protocols for consistent extraction of high quality DNA and
423 RNA from challenging plant and fungal samples for high-throughput SNP genotyping and sequencing
424 applications’, *PloS one*, 13(10), p. e0206085.
- 425 Jedlicka, P., Lexa, M. and Kejnovsky, E. (2020) ‘What Can Long Terminal Repeats Tell Us About the Age of
426 LTR Retrotransposons, Gene Conversion and Ectopic Recombination?’, *Frontiers in plant science*, 11, p. 644.
- 427 Jones, Philip, David Binns, Hsin-Yu Chang, Matthew Fraser, Weizhong Li, Craig McAnulla, Hamish
428 McWilliam, et al. (2014). “InterProScan 5: Genome-Scale Protein Function Classification.” *Bioinformatics* 30
429 (9): 1236–40.
- 430 Kebrom, T.H., McKinley, B. and Mullet, J.E. (2017) ‘Dynamics of gene expression during development and
431 expansion of vegetative stem internodes of bioenergy sorghum’, *Biotechnology for biofuels*, 10, p. 159.
- 432 Kolmogorov, M. *et al.* (2019) ‘Assembly of long, error-prone reads using repeat graphs’, *Nature*
433 *biotechnology*, 37(5), pp. 540–546.
- 434 Koren, S. *et al.* (2017) ‘Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and
435 repeat separation’, *Genome research*, 27(5), pp. 722–736.
- 436 Kumar, N. *et al.* (2022) ‘Registration of the sorghum carbon-partitioning nested association mapping (CP-
437 NAM) population’, *Journal of plant registrations* [Preprint]. doi:10.1002/plr2.20229.
- 438 Lexa, M. *et al.* (2020) ‘TE-greedy-nester: structure-based detection of LTR retrotransposons and their
439 nesting’, *Bioinformatics*, 36(20), pp. 4991–4999.
- 440 Li, H. *et al.* (2009) ‘The Sequence Alignment/Map format and SAMtools’, *Bioinformatics*, 25(16), pp. 2078–
441 2079.
- 442 Li, H. (2018) ‘Minimap2: pairwise alignment for nucleotide sequences’, *Bioinformatics*, 34(18), pp. 3094–
443 3100.
- 444 Li, J. *et al.* (2021) ‘Cotton pan-genome retrieves the lost sequences and genes during domestication and
445 selection’, *Genome biology*, 22(1), p. 119.
- 446 Li, X. *et al.* (2018) ‘Genomic and environmental determinants and their interplay underlying phenotypic
447 plasticity’, *Proceedings of the National Academy of Sciences of the United States of America*, 115(26), pp.
448 6679–6684.
- 449 Ma, J. and Bennetzen, J.L. (2004) ‘Rapid recent growth and divergence of rice nuclear genomes’, *Proceedings*
450 *of the National Academy of Sciences of the United States of America*, 101(34), pp. 12404–12410.
- 451 Marçais, G. *et al.* (2018) ‘MUMmer4: A fast and versatile genome alignment system’, *PLoS computational*
452 *biology*, 14(1), p. e1005944.
- 453 McCormick, Ryan F., Sandra K. Truong, Avinash Sreedasyam, Jerry Jenkins, Shengqiang Shu, David Sims,
454 Megan Kennedy, et al. (2018). “The Sorghum Bicolor Reference Genome: Improved Assembly, Gene
455 Annotations, a Transcriptome Atlas, and Signatures of Genome Organization.” *The Plant Journal: For Cell*
456 *and Molecular Biology* 93 (2): 338–54.
- 457 Multani, D.S. *et al.* (2003) ‘Loss of an MDR transporter in compact stalks of maize br2 and sorghum dw3
458 mutants’, *Science*, 302(5642), pp. 81–84.
- 459 Olson, Andrew J., and Doreen Ware. (2020). “Ranked Choice Voting for Representative Transcripts with
460 TRaCE.” *Cold Spring Harbor Laboratory*. <https://doi.org/10.1101/2020.12.15.422742>.

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

- 461 Ortiz, D., Hu, J. and Salas Fernandez, M.G. (2017) ‘Genetic architecture of photosynthesis in Sorghum bicolor
462 under non-stress and cold stress conditions’, *Journal of experimental botany*, 68(16), pp. 4545–4557.
- 463 Ou, S. *et al.* (2019) ‘Benchmarking transposable element annotation methods for creation of a streamlined,
464 comprehensive pipeline’, *Genome biology*, 20(1), p. 275.
- 465 Paterson, A.H. *et al.* (2009) ‘The Sorghum bicolor genome and the diversification of grasses’, *Nature*,
466 457(7229), pp. 551–556.
- 467 Quinlan, A.R. and Hall, I.M. (2010) ‘BEDTools: a flexible suite of utilities for comparing genomic features’,
468 *Bioinformatics*, 26(6), pp. 841–842.
- 469 Ruperao, Pradeep, Nepolean Thirunavukkarasu, Prasad Gandham, Sivasubramani Selvanayagam, Mahalingam
470 Govindaraj, Baloua Nebie, Eric Manyasa, *et al.* (2021). “Sorghum Pan-Genome Explores the Functional
471 Utility for Genomic-Assisted Breeding to Accelerate the Genetic Gain.” *Frontiers in Plant Science* 12 (June):
472 666342.
- 473 SanMiguel, P. *et al.* (1998) ‘The paleontology of intergene retrotransposons of maize’, *Nature genetics*, 20(1),
474 pp. 43–45.
- 475 Sayols, S. (2020) ‘rrvgo: a Bioconductor package to reduce and visualize Gene Ontology terms. 2020’.
- 476 Schliep, K.P. (2011) ‘phangorn: phylogenetic analysis in R’, *Bioinformatics*, 27(4), pp. 592–593.
- 477 Shumate, Alaina, and Steven L. Salzberg. (2021). “Liftoff: Accurate Mapping of Gene Annotations.”
478 *Bioinformatics* 37 (12): 1639–43.
- 479 Songsomboon, K. *et al.* (2021) ‘Genomic patterns of structural variation among diverse genotypes of Sorghum
480 bicolor and a potential role for deletions in local adaptation’, *G3* [Preprint]. doi:10.1093/g3journal/jkab154.
- 481 Stabenau, Arne, Graham McVicker, Craig Melsopp, Glenn Proctor, Michele Clamp, and Ewan Birney. (2004).
482 “The Ensembl Core Software Libraries.” *Genome Research* 14 (5): 929–33.
- 483 Tao, Yongfu, Hong Luo, Jiabao Xu, Alan Cruickshank, Xianrong Zhao, Fei Teng, Adrian Hathorn, *et al.*
484 (2021). “Extensive Variation within the Pan-Genome of Cultivated and Wild Sorghum.” *Nature Plants* 7 (6):
485 766–73.
- 486 Thompson, J.D., Higgins, D.G. and Gibson, T.J. (1994) ‘CLUSTAL W: improving the sensitivity of
487 progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and
488 weight matrix choice’, *Nucleic acids research*, 22(22), pp. 4673–4680.
- 489 Vaser, R. *et al.* (2016) ‘SIFT missense predictions for genomes’, *Nature protocols*, 11(1), pp. 1–9.
- 490 Vilella, Albert J., Jessica Severin, Abel Ureta-Vidal, Li Heng, Richard Durbin, and Ewan Birney. (2009).
491 “EnsemblCompara GeneTrees: Complete, Duplication-Aware Phylogenetic Trees in Vertebrates.” *Genome*
492 *Research* 19 (2): 327–35.
- 493 Waclawovsky, A.J. *et al.* (2010) ‘Sugarcane for bioenergy production: an assessment of yield and regulation
494 of sucrose content’, *Plant biotechnology journal*, 8(3), pp. 263–276.
- 495 Waititu, Joram Kiriga, Chunyi Zhang, Jun Liu, and Huan Wang. (2020). “Plant Non-Coding RNAs: Origin,
496 Biogenesis, Mode of Action and Their Roles in Abiotic Stress.” *International Journal of Molecular Sciences*
497 21 (21). <https://doi.org/10.3390/ijms21218401>.
- 498 Wang, Bo, Yinping Jiao, Kapeel Chougule, Andrew Olson, Jian Huang, Victor Llaca, Kevin Fengler, *et al.*

Ten new reference-quality genome assemblies for diverse bioenergy sorghum genotypes

- 499 (2021). “Pan-Genome Analysis in Sorghum Highlights the Extent of Genomic Variation and Sugarcane Aphid
500 Resistance Genes.” *bioRxiv*. <https://doi.org/10.1101/2021.01.03.424980>.
- 501 Wang, Y. *et al.* (2017) ‘Identification of tRNA nucleoside modification genes critical for stress response and
502 development in rice and Arabidopsis’, *BMC plant biology*, 17(1), p. 261.
- 503 Wayne Smith, C. and Frederiksen, R.A. (2000) *Sorghum: Origin, History, Technology, and Production*. John
504 Wiley & Sons.
- 505 Wu, X. *et al.* (2010) ‘Features of sweet sorghum juice and their performance in ethanol fermentation’,
506 *Industrial crops and products*, 31(1), pp. 164–170.
- 507 Wu, Y. *et al.* (2019) ‘Allelochemicals targeted to balance competing selections in African agroecosystems’,
508 *Nature plants*, 5(12), pp. 1229–1236.
- 509 Zhang, B. *et al.* (2019) ‘The poplar pangenome provides insights into the evolutionary history of the genus’,
510 *Communications biology*, 2, p. 215.
- 511 Zhang, L.-M. *et al.* (2018) ‘Sweet Sorghum Originated through Selection of Dry, a Plant-Specific NAC
512 Transcription Factor Gene’, *The Plant cell*, 30(10), pp. 2286–2307.
- 513 Zhou, Y. *et al.* (2020) ‘A platinum standard pan-genome resource that represents the population structure of
514 Asian rice’, *Scientific data*, 7(1), p. 113.
- 515 Zhou, Y. *et al.* (2022) ‘Graph pangenome captures missing heritability and empowers tomato breeding’,
516 *Nature*, 606(7914), pp. 527–534.

517 Data Availability Statement

518 Assembled Genomes are publicly available on <https://www.sorghumbase.org/>. Gene data hosted at
519 https://ftp.sorghumbase.org/Voelker_et_al_2022/. Raw data and genome assemblies are available at
520 the European Nucleotide Archive under the project ID: PRJEB55613

521