

GigaScience

Assembly of the 373K gene space of the polyploid sugarcane genome reveals reservoirs of functional diversity in the world's leading biomass crop

--Manuscript Draft--

Manuscript Number:	GIGA-D-19-00013	
Full Title:	Assembly of the 373K gene space of the polyploid sugarcane genome reveals reservoirs of functional diversity in the world's leading biomass crop	
Article Type:	Research	
Funding Information:	FAPESP (2012/51062-3)	Professor Glaucia Mendes Souza
	FAPESP (2008/52146-0)	Professor Glaucia Mendes Souza
	FAPESP (2014/50921-8)	Professor Glaucia Mendes Souza
	FAPESP (2008/52074-0)	Not applicable
	FAPESP (2011/50761-2)	Not applicable
	National Science Foundation (DBI-1350041)	Not applicable
	CNPq (304360/2014-7)	Professor Glaucia Mendes Souza
	CNPq (308197/2010-0)	Not applicable
	FAPESP (2015/22993-7)	Not applicable
	FAPESP (2013/18322-4)	Not applicable
	FAPESP (2015/15346-5)	Not applicable
	CNPq (159094/2014-3)	Not applicable
	FAPESP (2017/02270-6)	Not applicable
	CAPES (DS-1454337)	Not applicable
	FAPESP (2013/23048-9)	Not applicable
	FAPESP (2016/06917-1)	Not applicable
	FAPESP (2013/07467-1)	Not applicable
	FAPESP (2017/02842-0)	Not applicable
	CNPq (309566/2015-0)	Not applicable
	National Science Foundation (IOS/0115903)	Not applicable
	National Institutes of Health (R01-HG006677)	Not applicable
Abstract:	<p>Background Sugarcane cultivars are polyploid interspecific hybrids of giant genomes, typically with 10-13 sets of chromosomes from two <i>Saccharum</i> species. The ploidy, hybridity and size of the genome, estimated to have in excess of 10 Gb, pose a great challenge for sequencing.</p> <p>Results Here we present a gene-space assembly of SP80-3280, including 373,869 putative genes and their potential regulatory regions. Their alignment to single copy genes of diploid grasses indicates that we could resolve 2-6 (up to 15) gene copies</p>	

	<p>(homo/homeolog) that are 99.1% identical within their coding sequences. Dissimilarities increase in their regulatory regions and gene promoter analysis shows differences in regulatory elements within gene families and are species-specific expressed. We exemplify these differences for sucrose synthase (SuSy) and phenylalanine ammonia-lyase (PAL), two gene families central to carbon partitioning. SP80-3280 have particular regulatory elements involved in sucrose synthesis not found in the ancestor <i>S. spontaneum</i>. PAL regulatory elements are found in co-expressed genes related to fiber synthesis within gene networks defined during plant growth and maturation. Comparison to sorghum reveals predominantly biallelic variations in sugarcane, consistent with the formation of two 'subgenomes' after their divergence ca. 3.8~4.6 MYA and reveals SNVs that may underlie their differences.</p> <p>Conclusions</p> <p>This gene-copy resolved assembly represents a large step towards a whole genome assembly of a commercial sugarcane cultivar providing a large diversity of genes and homo(eo)logs useful for improving biomass and food production.</p>
Corresponding Author:	<p>Glauca Mendes Souza, Ph.D Universidade de São Paulo Sao Paulo, SP BRAZIL</p>
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	<p>Universidade de São Paulo</p>
Corresponding Author's Secondary Institution:	
First Author:	<p>Glauca Mendes Souza, Ph.D</p>
First Author Secondary Information:	
Order of Authors:	<p>Glauca Mendes Souza, Ph.D</p> <p>Marie-Anne Van Sluys, Ph.D</p> <p>Carolina Gimiliani Lembke, Ph.D</p> <p>Hayan Lee, Ph.D</p> <p>Gabriel Rodrigues Alves Margarido, Ph.D</p> <p>Carlos Takeshi Hotta, Ph.D</p> <p>Jonas Weissmann Gaiarsa, Ph.D</p> <p>Augusto Lima Diniz, Ph.D</p> <p>Mauro de Medeiros Oliveira, Ph.D</p> <p>Sávio de Siqueira Ferreira, Ph.D</p> <p>Milton Yutaka Nishiyama-Jr, Ph.D</p> <p>Felipe ten Caten, Ph.D</p> <p>Geovani Tolfo Ragagnin, MSc</p> <p>Pablo de Moraes Andrade, Ph.D</p> <p>Robson Francisco de Souza, Ph.D</p> <p>Gianluca Gonçalves Nicastro, Ph.D</p> <p>Ravi Pandya, BS.c</p> <p>Changsoo Kim, Ph.D</p> <p>Hui Guo, Ph.D</p> <p>Alan Mitchell Durham, Ph.D</p> <p>Monalisa Sampaio Carneiro, Ph.D</p> <p>Jisen Zhang, Ph.D</p> <p>Qing Zhang, Ph.D</p>

	Qing Zhang, Ph.D
	Ray Ming, Ph.D
	Michael Schatz, Ph.D
	Bob Davidson
	Andrew Paterson, Ph.D
	David Heckerman, Ph.D
Order of Authors Secondary Information:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the</p>	Yes

conclusions of the paper rely must be either included in your submission or deposited in [publicly available repositories](#) (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

[Click here to view linked References](#)

1 Assembly of the 373K gene space of the polyploid sugarcane genome reveals reservoirs of 1 2 functional diversity in the world's leading biomass crop

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

Full name	Institutional address	e-mail
Glaucia Mendes Souza*	1	glmsouza@iq.usp.br
Marie-Anne Van Sluys*	2	mavsluys@usp.br
Carolina Gimiliani Lembke	1	carolina.lembke@gmail.com
Hayan Lee	3,4	hayan.lee@stanford.edu
Gabriel Rodrigues Alves Margarido	5	gramarga@usp.br
Carlos Takeshi Hotta	1	hotta@iq.usp.br
Jonas Weissmann Gaiarsa	2	jonaswg@gmail.com
Augusto Lima Diniz	1	augustold@usp.br
Mauro de Medeiros Oliveira	1	mauromedeiros@usp.br
Sávio de Siqueira Ferreira	1,2	saviobqi@gmail.com
Milton Yutaka Nishiyama-Jr	1,6	yutakajr@gmail.com
Felipe ten Caten	1	ftencaten@gmail.com
Geovani Tolfo Ragagnin	2	geovaniragagnin@gmail.com
Pablo de Morais Andrade	1	pablo.andrade@gmail.com
Robson Francisco de Souza	7	rfsouza@usp.br
Gianluca Gonçalves Nicastro	7	nicastro@iq.usp.br
Ravi Pandya	8	ravip@microsoft.com,
Changsoo Kim	9,10	changsookim@cnu.ac.kr
Hui Guo	9	huiguo7@gmail.com
Alan Mitchell Durham	11	aland@usp.br
Monalisa Sampaio Carneiro	12	monalisa@ufscar.br
Jisen Zhang	13	zjisen@126.com
Xingtang Zhang	13	tanger_009@163.com
Qing Zhang	13	zhangqing970@126.com
Ray Ming	13,14	rayming@illinois.edu
Michael C. Schatz	3,15	michael.schatz@gmail.com
Bob Davidson	8	bob.davidson@microsoft.com
Andrew Paterson	9	paterson@uga.edu
David Heckerman	8	heckerma@hotmail.com

1 – Departamento de Bioquímica, Instituto de Química, Universidade de São Paulo, Av. Prof. Lineu Prestes,

748, São Paulo, SP 05508-000, Brazil

2 – Departamento de Botânica, Instituto de Biociências, Universidade de São Paulo, Rua do Matão, 277, São

Paulo, SP 05508-090, Brazil

3 – Cold Spring Harbor Laboratory, One Bungtown Road, Koch Building #1119, Cold Spring Harbor, NY

11724, United States of America

- 11 4 – Department of Energy Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, CA CA 94598, United
12 States of America
2
313 5 – Departamento de Genética, Escola Superior de Agricultura Luiz de Queiroz, Universidade de São Paulo,
4
514 Avenida Pádua Dias, 11, Piracicaba, SP 13418-900, Brazil
6
715 6 – Laboratório Especial de Toxinologia Aplicada, Instituto Butantan, Av. Vital Brasil, 1500, São Paulo, SP
8
916 05503-900, Brazil
10
1117 7 – Departamento de Microbiologia, Instituto de Ciências Biomédicas, Universidade de São Paulo, Av.
13
1418 Professor Lineu Prestes, 1734, São Paulo, SP 05508-900, Brazil
15
1619 8 – Microsoft Research, One Microsoft Way, Redmond, WA 98052, United States of America9 – Plant
17
1820 Genome Mapping Laboratory, University of Georgia, 120 Green Street, Athens, GA 30602-7223, United
19
2021 States of America
21
2222 10 – Department of Crop Science, Chungnam National University, 99 Daehak Ro Yuseong Gu, Deajeon,
24
2523 34134, South Korea
26
2724 11 – Departamento de Ciências da Computação, Instituto de Matemática e Estatística, Universidade de São
28
3025 Paulo, Rua do Matão, 1010, São Paulo, SP 05508-090, Brazil
31
3226 12 - Departamento de Biotecnologia e Produção Vegetal e Animal, Centro de Ciências Agrárias, Universidade
33
3427 Federal de São Carlos, Rodovia Washington Luis km 235, Araras, SP 13.565-905, Brazil
35
3628 13 – FAFU and UIUC-SIB Joint Center for Genomics and Biotechnology, Fujian Agriculture and Forestry
37
3829 University, Shangxiadian Road, Fuzhou 350002, Fujian, China
40
4130 14 - Department of Plant Biology, University of Illinois at Urbana-Champaign, 201 W. Gregory Dr. Urbana,
42
4331 Urbana, Illinois 61801, USA
44
4532 15 – Departments of Computer Science and Biology, Johns Hopkins University, 3400 North Charles Street,
46
4733 Baltimore, MD 21218-2608, United States of America
48
49
5034
51

5235 *These authors contributed equally to this work and are co-corresponding authors: glmsouza@iq.usp.br and
53
5436 mavsluys@usp.br
55

56
5737
58
5938
60
6139
62
63
64
65

40 **ABSTRACT**

41

2

342 **Background**

4

543 Sugarcane cultivars are polyploid interspecific hybrids of giant genomes, typically with 10-13 sets of
6
744 chromosomes from two *Saccharum* species. The ploidy, hybridity and size of the genome, estimated to have
8
9
1045 in excess of 10 Gb, pose a great challenge for sequencing.

11

1246 **Results**

13

1447 Here we present a gene-space assembly of SP80-3280, including 373,869 putative genes and their potential
15
1648 regulatory regions. Their alignment to single copy genes of diploid grasses indicates that we could resolve 2-
17
18
1949 6 (up to 15) gene copies (homo/homeolog) that are 99.1% identical within their coding sequences.
20
2150 Dissimilarities increase in their regulatory regions and gene promoter analysis shows differences in regulatory
22
2351 elements within gene families and are species-specific expressed. We exemplify these differences for sucrose
24
2552 synthase (SuSy) and phenylalanine ammonia-lyase (PAL), two gene families central to carbon partitioning.
26
2753 SP80-3280 have particular regulatory elements involved in sucrose synthesis not found in the ancestor *S.*
28
29
3054 *spontaneum*. PAL regulatory elements are found in co-expressed genes related to fiber synthesis within gene
31
3255 networks defined during plant growth and maturation. Comparison to sorghum reveals predominantly biallelic
33
3456 variations in sugarcane, consistent with the formation of two ‘subgenomes’ after their divergence ca. 3.8~4.6
35
36
3757 MYA and reveals SNVs that may underlie their differences.

38

3958 **Conclusions**

40

4159 This gene-copy resolved assembly represents a large step towards a whole genome assembly of a commercial
42
4360 sugarcane cultivar providing a large diversity of genes and homo(eo)logs useful for improving biomass and
44
45
4661 food production.

47

4862

49

5063 **Keywords:** Allele; Bioenergy; Biomass; Genome; Polyploid

51

5264

53

5465

55

5666

57

5867

60

6168

62

63

64

65

69 BACKGROUND

170 Sugarcane is the world's most cultivated crop in tonnage (more than rice, maize and wheat) [1], and is
2
371 considered the most sustainable of energy crops [2] with high potential to mitigate climate change without
4
572 affecting food security [3]. Already produced in over 100 countries, high productivity of sugar, bioethanol and
6
773 bioelectricity [4] make it a highly expandable green alternative to petroleum [5–7]. The International Energy
8
974 Agency projects a 150 EJ (17% of energy demand) contribution of bioenergy by 2060, delivering 18% of the
11
1275 emission reductions needed to achieve the 2DS (2°C Scenario). Sugarcane bioenergy production by 2045 could
13
1476 displace up to 13.7% of crude oil consumption and 5.6% of the world's CO₂ emissions relative to 2014. This
15
1677 can be achieved without using forest preservation areas or land necessary for food production systems.
17
1878 Additionally, the myriad of products that can derive from sugarcane biomass [8] further enhance opportunities
19
20
2179 for sugarcane in a portfolio of technologies needed to transition to a low carbon 'bioeconomy'.

22
2380 Opportunities to accelerate breeding progress and enrich knowledge of the fundamental biology of this
24
2581 important plant motivate efforts to produce a high-quality reference genome, a challenge that is unusually
26
2782 complex. Unlike wheat cultivated species known to be either tetraploid (AABB) or hexaploid (AABBDD), the
28
29
3083 *Saccharum* (sugarcane) genus is considered to be a species complex. A recent study [9] proposed independent
31
3284 polyploidization events within *Saccharum* after divergence from the last ancestor shared with *Sorghum*,
33
3485 superimposed upon an additional whole genome duplication since the diversification of grasses. As a
35
3686 consequence, the sugarcane genome is redundant and harbors genes in multiple functional copies. Adding
37
38
3987 further complexity, sugarcane cultivars are polyploid interspecific hybrids, typically with 10-13 sets of their
40
4188 10 basic chromosomes, 80-85% from *Saccharum officinarum* (2n=80), which is known for its sweetness, 10-
42
4389 15% from *S. spontaneum* (2n=40-128) known for its robustness, and ~5% with recombined chromosomes
44
4590 between those two progenitors [10,11]. The ploidy, hybridity and sheer size of the genome, estimated to have
46
47
4891 in excess of 10 Gb, pose a great challenge for sequencing [12]. Recently released sequences of the modern
49
5092 cultivar R570 yielded a mosaic monoploid reference (382 Mb single tiling path) [13] and a *S. spontaneum*
51
5293 AP85-441 haploid assembly (3.13 Gb) [14].

53
5494 Worldwide sugarcane yield (~84 ton/ha) is currently only ~20% of the theoretical potential yield (~381
55
5695 ton/ha), spurring great interest in conventional or molecular breeding approaches to improve yield. However,
57
58
5996 progress by conventional breeding towards closing the gap between current and potential yield has been slow
60
6197 with yield gains in the order of 1.0–1.5% a year [15]. Sugarcane commercial cultivars distribute roughly one

98 third of their carbon into sucrose and two thirds into tops and stems which, due to high lignin content, are
99 burned to fuel boilers, contributing to the favorable energy balance of industrial processes [16]. As sugarcane
100 can accumulate large amounts of sucrose in its stems, up to ~650 mM [17], it is important to study sucrose
101 metabolism and the key players in its regulation. Also, of interest is the revealing of regulators of cell wall
102 biosynthesis. Altering these pathways may help shift carbon partitioning from sucrose storage to biomass
103 accumulation, rich in fiber content, mostly composed of secondary cell walls formed by cellulose,
104 hemicellulose and lignin [18]. The latter compound is a hydrophobic polymer that provides strength and
105 rigidity to the plant, but also is responsible for cell wall recalcitrance, which is the natural plant resistance to
106 hydrolytic attacks that hampers cellulosic ethanol production [19].

109 RESULTS

111 The SP80-3280 assembly reveals a gene space of 373,869 genes

112 Here, we report a representative gene space assembly of the genome sequence of SP80-3280 (GenBank
113 accession number QPEU01000000), the cultivar used in Brazilian breeding programs with the largest
114 collection of transcriptomic data available [20]. On average, 6 sugarcane haplotypes, putatively homo(eo)logs,
115 could be resolved from 4.26 Gb of assembled data and 373,869 putative genes and promoter regions. This is
116 the first release of a resolved gene assembly of such a giant hybrid polyploid genome and their potential
117 regulatory regions.

118 The assembly was constructed using 26 libraries sequenced using Illumina Synthetic Long-Read
119 technology, obtaining 19 Gb, ~19x haploid genome coverage (~1.9X genome coverage) with >99% of bases
120 having >99% accuracy (**Additional file 1: Fig. S1**). The final assembly includes 450,609 contigs (unitigs +
121 singletons), with average length of 9,452 bp and NG50 of 41,394 bp (**Table 1**), adding over 3Gb of sequence
122 not previously reported (**Additional file 2: Table S1**) [21]. The gene space described here is organized in the
123 SUCEST-FUN public database (<http://sucest-fun.org/wsapp/>).

124 Several indicators support the comprehensiveness of the SP80-3280 gene space: (i) among 39,441
125 sorghum transcripts, 39,207 (99.4%) matched the assembly, at least partially; of these, 71.1% matched at least
126 one sugarcane contig with 90% or higher coverage (**Additional file 1: Fig. S2**); (ii) the assembly completely

127 covers 217 (87.5%) of the 248 ultra-conserved Core Eukaryotic Genes Mapping Approach (CEGMA) [22]
128 proteins, and partly covers 18 (7.3%), with only 13 (5.2%) not detected (**Additional file 2: Table S2**); (iii)
129 among 1,440 genes in the Benchmarking Universal Single-Copy Orthologs (BUSCO) [23] Plantae lineage, the
130 assembly completely covers 1,309 (90.9%) and partially covers 53 (3.7%) (**Additional file 2: Table S3**). By
131 including tBLASTn of the 78 (5.4%) missing Plantae lineage BUSCO genes, only 8 (0.5%) are absent; (iv)
132 assembled chloroplast (MG969494) and mitochondrial (MG969495 and MG969496) genomes were over 99%
133 similar (at gene level) to published *Saccharum* genomes [24–26]; and (v) 94.9% of 134,840 SP80-3280
134 expressed sequence tags (ESTs) match the assembled gene-space sequence.

135 The assembly revealed 373,869 putative genes with 374,774 transcripts (**Table 1**), far more than the
136 72,269 unigenes inferred from six sugarcane genotypes [27]; 85,151 transcripts of sugarcane genotypes with
137 contrasting lignin contents [28]; and 195,765 transcripts inferred from *de novo* assembly of ORFeomes from
138 *S. officinarum*, *S. spontaneum* and SP80-3280 [29]. The number of genes, high quality of alignments, and the
139 following analysis indicates that the assembly provides a high-quality resolution of homo(eo)logs genes.

140 Among the predicted transcripts, 302,627 (80.7%) aligned to a Uniref50 protein [30], and 195,651 were
141 annotated with 10,362 GO terms [31] (**Additional file 1: Fig. S3**). Our previously published SP80-3280
142 ORFeome was reassembled using the genome as a reference, revealing 269,050 genes and 275,807 transcripts
143 from leaves, immature and intermediate internodes (**Additional file 2: Table S4**). Further, a set of 134,840
144 SP80-3280 ESTs from a Sugarcane EST Project – SUCEST [20] – were mapped to assembled contigs and
145 compared to predicted genes, in order to further estimate the completeness of the predicted gene space. A total
146 of 127,940 ESTs (94.9%) have at least one match in the assembly, which is in accordance with similar analysis
147 of other plant genomes [32], and only 7.6% of aligned ESTs (9,772) do not correspond with predicted genes.
148 This result resembles the BUSCO results, for which only 5.4% of conserved genes could not be identified in
149 the assembly. Performing the same approach using a set of 43,141 sugarcane assembled sequences (SAS) –
150 derived from an EST project [20] – produces similar results regarding the overall alignment rate, with 93% of
151 the sequences (40,147) matching at least one location in contigs. It is important to note that the SUCEST
152 database included ESTs from 10 sugarcane cultivars including SP80-3280 which may explain the reduced
153 correspondence with genic regions, with 16% of aligned SAS having no correspondence to predicted genes of
154 SP80-3280.

155 To verify how the assembled genes reflected the expected content of homo(eo)logs genes, the gene content
156 was compared to those of other grasses. Single-copy genes from diploid grasses (sorghum, rice and
2
157 *Brachypodium*) are present in up to 15 copies in sugarcane, mostly with 2-6 copies (total of 1,592 coding
4
158 sequences (CDS) in sugarcane) (**Fig. 1A**). Dissimilarities among gene copies increase from the coding region
6
159 to the promoter region, with median divergence of 0.90% between CDS, 1.03% for the 100 nucleotides (nt)
8
160 upstream, 4.47% for 500 nt and 7.50% for 1,000 nt (**Fig. 1B**). Frame-preserving INDELS are more abundant
10
161 than frameshifts (**Fig. 1C**) and short frameshift INDELS were relatively less frequent in the sugarcane exons
13
162 than in sorghum [33].
15

16
163
17

164 **Differential homo(eo)logs expression is observed**

20
2165 The SP80-3280 gene series that correspond to diploid grass single-copy genes showed expression of sense
22
2166 copies for multiple homo(eo)logs (**Fig. 2A**), with very few copies transcribed in antisense orientation (**Fig.**
24
2167 **2B**) based on alignment with the SP80-3280 cDNA reads [29] from leaves, immature and intermediate
26
2168 internodes. For some genes, not all copies are expressed in SP80-3280 (**Fig. 2A, Additional file 1: Fig. S4**).
29
3169 The number of expressed homo(eo)logs is different across the three tissues (**Additional file 1: Fig. S4A**). This
31
3170 difference among copies is consistent with the divergence of upstream regions (putatively gene promoter
33
3171 regions). The increase in the number of expressed copies is not accompanied by an increase in the level of
35
3172 expression (**Additional file 1: Fig. S4B**).
38

3973 As an example of the complexities in data mining of such a complex gene-space for future reference, we
40
4174 offer an example using some well-known genes involved in biomass production.
42

43
44
45

4676 **Gene family analysis of SuSy and PAL shows differences in their regulatory regions in SP80-3280 and** 47 4877 ***S. spontaneum***

49
5178 Sucrose Synthases (SuSy) catalyze the reversible breakdown of sucrose into UDP-glucose and fructose in
51
5179 carbon partitioning [34]. In agreement with previous work on sugarcane progenitors [35] (*S. officinarum*, *S.*
53
5180 *robustum* and *S. spontaneum*), phylogenetic analysis of 44 ScSuSy members identified in the SP80-3280
56
5181 assembly supports that this hybrid has 5 SuSy genes (hereafter ScSuSy1-5) in three clades: I (ScSuSy1 and 2),
58
5182 II (ScSuSy3 and 5) and III (ScSuSy4) (**Fig. 3A**). Sorghum shares these 5 SuSy genes, indicating that they
60
61
62
63
64
65

183 evolved before the sugarcane/sorghum divergence. RNA-Seq data from leaves and internodes of SP80-3280
184 [29] shows expression of 36 of the 44 ScSuSy members, suggesting ScSuSy1-2 (clade I) and ScSuSy5 might
2
185 control carbon flux from source to biomass conversion in stems, as they show higher expression in internodes
4
186 than in leaves (**Fig. 3C**).

187 SuSy produces the substrate for cellulose biosynthesis (UDP-glucose) and is commonly associated with
9
188 cell wall and cellulose synthesis [36,37]. In view of the myriad of possibilities to convert lignocellulosic
11
189 compounds into chemicals and fuels, defining phenylpropanoid biosynthesis pathway members in sugarcane
13
190 is of great interest. Phenylalanine ammonia-lyase (PAL), the first enzyme in phenylpropanoid biosynthesis
15
191 [38–40], is correlated with lignin content [38–41], a major component of plant cell walls [18], and is responsive
18
192 to the ethylene-releasing ripener (ethephon) in both leaf and internode [42].

193 Mapping of predicted proteins from SP80-3280 against the SUCEST-FUN Cell Wall Catalogue [41] (731
22
194 transcripts of 20 protein categories) identified 3,054 similar proteins (**Additional file 2: Table S5**), including
25
195 47 PAL copies. Phylogenetic analysis together with sorghum, *S. spontaneum* and mosaic monoploid R570
27
196 PAL sequences reveals 5 clusters (**Fig. 3B**), each containing at least one representative with a sorghum
29
197 ortholog. *S. spontaneum* has 33 putative PAL genes, somewhat more than expected considering that the
31
198 sequenced genotype is a tetraploid. The higher number may be due to expansion of PAL members in clade I
34
199 that occurred also for sorghum and the sugarcane hybrid genomes of R570 and SP80-3280. Clade V has a
36
200 higher number of SP80-3280 PAL members and all except one (ID 37780.4) showed expression evidence (**Fig.**
38
201 **3D**). In addition, the CCR (Cinnamoyl-CoA reductase), COMT (Caffeic acid 3-O-methyltransferase) and 4CL
40
202 (4-coumarate-CoA ligase) gene families, also related to phenylpropanoid biosynthesis, have much higher
43
203 numbers of genes (620, 453 and 375, respectively) in sugarcane than sorghum [43] (44, 41 and 15,
45
204 respectively). This is another challenge and opportunity for future functional characterization (**Additional file**
47
205 **2: Table S6**).

206 The sheer number of sugarcane genes found so far, the large size of multi-gene families and the finding
52
207 that homo(eo)logs are differentially expressed point to a very complex role of regulation in the determination
54
208 of phenotypic differences. Consistent with the gene copy-richness of sugarcane, we inferred 15,737
56
209 transcription factors (TFs) from 57 families (**Additional file 2: Table S7**), versus ~2,000 previously estimated
58
210 [44]. The classification of core promoters and identification of Transcription Factor Binding Sites (TFBSs) in
61
62
63
64
65

211 proximal promoters was performed *in silico* and the percentage of core promoter regions with a TATA-box
212 element was 47.72% and 12.76% for SuSy and PAL genes, respectively. The TFBS identification pointed to a
213 wealth of regulatory elements differentially distributed among members of the same gene family, i.e. SuSy
214 and PAL (**Fig. 3C and D and Additional file 2: Table S8**). In addition, using gene expression data of SP80-
215 3280 plants grown in field conditions for 13 months, we have found evidence of a co-expression module,
216 enriched for phenylpropanoid and lignin biosynthesis gene ontology terms (**Additional file 1: Fig. S5A**). This
217 module comprises 116 transcripts, including one PAL (**Additional file 1: Fig. S5B**), whose expression is
218 higher in internodes 5 and 9, than in leaves and immature internode (**Additional file 1: Fig. S5C**). It was
219 possible to identify the TFBSs, predicted as putative regulators of the PAL gene family (**Fig. 3D**) within the
220 upstream region of these co-expressed genes, suggesting that ABF, ERF, ZF-HD/C2H2, and ARF3
221 (**Additional file 1: Fig. S5D**) may also regulate other genes involved in lignin biosynthesis and metabolism.

222 The most significant motifs found for each gene family (SuSy and PAL) were mapped to the promoter
223 region of the remaining sequences from both SP80-3280 and R570 hybrids and *S. spontaneum* (**Additional**
224 **file 2: Table S8c and Table S9**). Interestingly, only ScSuSy2 and ScSuSy3 motifs mapped in all species,
225 suggesting that SP80-3280 hold particular regulatory elements involved in sucrose synthesis. Conversely,
226 SP80-3280 and *S. spontaneum* share all predicted motifs for PAL genes (**Additional file 2: Table S9**),
227 suggesting that this gene family may be derived from the *S. spontaneum* ancestor.

228 **Transposable element insertions may affect SuSy and PAL expression**

229 Fewer transposable elements (TE) were identified in SP80-3280 gene space than in the AP85-441 *S.*
230 *spontaneum* and mosaic monoploid R570 assembly, probably due to repetitive regions collapsing in the
231 assembly even with the use of long synthetic-read sequencing (**Additional file 1: Fig. S6, Additional file 2:**
232 **Table S10**). All previously described TE families are represented in the three genome assemblies, disclosing
233 few cultivar specific amplifications. The two modern cultivars have fewer TE counts than the *S. spontaneum*
234 progenitor in normalized monoploid genomes. LTR retrotransposons are large contributors to genome
235 composition at the chromosome assembly level. However, scMaximus (Copia) and scDel (Gypsy) LTR-
236 retrotransposon families are similarly represented in both gene-space and chromosome assemblies supporting
237 their presence in transcriptionally active regions [45]. We also note that scCACTA transposons are more
238 represented at the gene-space assembly than schAT while the scMutator family is similarly represented in both.

240 Functionally important TE insertions were identified in the ScSuSy gene family (**Fig. 3**). ScSuSy2
241 copies have a contrasting pattern, most *S. spontaneum* having TE insertions while most SP80-3280
242 homo(eo)logs do not – although SP80-3280 and *S. spontaneum* share one ancient insertion of schAT159 at
243 similar distances from the ATG. ScSuSy3 genes are polymorphic between species and within SP80-3280, with
244 6 copies having no TE and 5 in which different TEs may impact expression. In particular,
245 scga7_uti_cns_0020964:7575-17575 (-) harbors a full LTR at 280 bases from the ATG. Most ScSuSy4 copies
246 have no TE insertion but interestingly, as described for ScSuSy2, SP80-3280 (scga7_uti_cns_0226458:7638-
247 16073 (-)) and *S. spontaneum* (Chr1B:33406669-33416669 (-)) share one ancient schAT159 insertion. Finally,
248 ScSuSy1 has similar patterns of TE presence and absence in both genomes, and ScSuSy5 genes have no
249 insertions in the promoter regions of either *S. spontaneum* or SP80-3280. Furthermore, PAL genes from clade
250 I exhibit most of the copy variation and harbor TEs inserted near the promoter region. Only two copies from
251 SP80-3280 and *S. spontaneum* lack TE insertion in PALs from group I.

253 **Sugarcane and sorghum polymorphisms support recent allotetraploidy**

254 Despite a common foundation for evolving high sugar content with similar Susy genes (ScSuSy1-5),
255 sugarcane and closely related sorghum have taken different paths since sharing ancestry. We identified 4,750
256 natural SNP variations (SNVs) between sorghum and sugarcane gene regions, mostly bi-allelic (3,840
257 (80.8%)), but 6.2% tri-allelic (295) and 0.97% tetra-allelic (46) (**Fig. 4**). Further, 1,334 SNVs that differentiate
258 sugarcane from sorghum in 585 single copy genes include frameshifts, premature splicing, loss of stop codons
259 and translation initiation (**Additional file 1: Fig. S7, Additional file 2: Table S11**) in genes significantly
260 enriched in transcription, DNA-dependent cell organization and biogenesis in the nucleus and endoplasmic
261 reticulum (**Additional file 2: Table S12**) comprise a rich slate of candidates for causes of morphological and
262 physiological differences between these taxa.

264 **The gene-space contribution towards a chromosome level assembly of a sugarcane commercial hybrid**

265 Notwithstanding the fragmented nature of our assembly, we explored how it could contribute beyond the
266 gene space toward a whole genome assembly of the hybrid sugarcane genome. Previous analysis of grass
267 genomes revealed extensive conservation of gene order overlaid with a background of small-scale
268 chromosomal rearrangements and numerous localized gene deletions, insertions and duplications [46].

269 Recently published estimates of the levels of gene synteny between *Sorghum bicolor* and the sugarcane cultivar
270 R570 found that 83% of the genes are arranged co-linearly in the two genomes [13]. In our assembly of SP80-
271 3280, 79,094 contigs had at least two predicted genes and could therefore be used to compare the order of
272 genes in SP80-3280 to those of sorghum. To avoid the need to resolve multiple comparisons to duplicated
273 regions in the sorghum genome, we generated a sequence similarity-based clustering of all coding sequences
274 from both genomes and used the genes in clusters with only one sorghum gene as anchors to evaluate synteny
(**Additional file 1: Fig. S8**). We found that 9,319 SP80-3280 contigs had at least two synteny anchors and
276 85% (7,906) of these contigs were fully syntenic (**Additional file 1: Fig. S9A, B**), *i.e.* had all genes in the
277 same order and orientation in SP80-3280 contigs and the sorghum chromosomes (**Additional file 2: Table**
278 **S13**). To evaluate the effect of SP80-3280 assembly fragmentation on the distribution of contigs per syntenic
279 block length, we sampled 10,000 contigs from SP80-3280 and 10,000 chromosome fragments from both R570
280 and *S. spontaneum*, while preserving the overall distribution of contig lengths observed in SP80-3280. The
281 distributions converged to similar shapes, with most contigs and chromosome fragments harboring a single
282 syntenic block in all genomes (**Additional file 1: Fig. S9C**). While the number of syntenic blocks per contig
283 were identical for the two cultivars, a larger frequency of fully syntenic contigs was observed for SP80-3280,
284 suggesting that our assembly is enriched in genomic neighborhoods that are co-linear to sorghum, presumably
285 comprising euchromatin. While an increase in sequencing coverage would lead to improved estimates of co-
286 linearity, our results agree with widespread findings on the conservation of gene order among grass species
287 and support the conclusion that, albeit fragmented, our assembly does not contain an excess of chromosomal
288 rearrangements, as would be expected if there was a significant amount of chimeric contigs.

289 Finally, to allocate the gene space into potential physical groupings we aligned the SP80-3280
290 transposable element (TE) masked BWA-SW to chromosome level assemblies of the *S. spontaneum* tetraploid
291 AP85-441 genome [14] and the R570 [13] monoploid genome data. Multiple correspondence analysis (MCA)
292 with hierarchical clustering of the sequences enabled us to allocate the gene space contigs into 6 clusters, an
293 important contribution to future scaffolding efforts. From the total of 450,609 contig sequences, 418,471
294 (92,86%) produced a BWA-SW alignment against the *S. spontaneum* [14] and R570 [13] assemblies (**Fig. 5A**)
295 and protein alignment among these three species are consistent with MCA results (**Fig. 5B and C**). Contigs
296 were also mapped against a collection of 778 targeted sequenced BACs of which 347 are from SP80-3280 and
297 431 from R570. All BACs had a corresponding contig match against the assembly. This collection shows

298 centromeric regions and non-TE multigene families are the most covered (64x). An R gene locus (I2C-2) found
299 in cluster 3 of SP80-3280 and in chromosome 9 of R570, was verified for co-location with a Ca⁺-dependent
300 kinase, a *dog1* (delay of germination 1) and an aminotransferase. The co-location was confirmed in R570 and
301 SP80-3280 BACs showing up to eight copies of each gene (**Additional file 1: Fig. S10**).

304 DISCUSSION

305 This assembly presents 373,869 genes. The gene space described here represents a significant step in
306 understanding the haplotype origin of the hybrid genome. Approximately 12.25% of the SP80-3280 sequences
307 are of *S. spontaneum* origin [14], supporting previous studies [10,11]. The comparison against different sets of
308 genes (sorghum, CEGMA, BUSCO, mitochondrial and chloroplast) supported the comprehensiveness of the
309 gene space. The total of predicted genes (373,869) is around 10x, 14x and 13x higher than those for monoploid
310 genome assemblies of *S. spontaneum* [14], sugarcane R570 [13] and sorghum [48], respectively. This is in
311 agreement with the predicted 8 to 14 copies for *S. spontaneum*, depending on the cytotypes, and for modern
312 sugarcane varieties [49]. Genes that are single-copy in diploid grasses are present in up to 15 copies in the
313 SP80-3280 assembly and the sequence differences are present mainly in the upstream regulatory region. This
314 highlights the importance and complexity of studying homo(eo)logs expression in sugarcane and adds great
315 value to the development of molecular markers for breeding in gene promoter regions. The differences in gene
316 upstream sequences causes the differential expression among the copies and across the studied tissues. This
317 was also reported for the polyploids cotton [50] and wheat [51]. Expression differences among homo(eo)logs
318 in polyploid species may play a crucial role in increasing adaptability to environmental stresses (such as
319 salinity [52], heat and drought [53]) and in improving performance of new cultivars. These differences
320 highlight the importance of our assembly which discriminates homo(eo)logs, for example providing
321 information important for the selection of target sequences (genes or promoters) to produce transgenic
322 sugarcane plants. With the homo(eo)logs identified, one could discard a sequence that is not expressed or use
323 genome editing tools to modify a target sequence to increase its expression. It is also possible to identify the
324 progenitor contributing a homo(eo)log (e.g., *S. spontaneum*, *S. officinarum* or a parent in a cross) and select
325 the homo(eo)log from the progenitor that has the phenotype of interest.

326 We also show how the data can be useful for gene promoter analysis. Expansion of SuSy genes might by
327 selected for fiber development in cotton [54]. Different members of the SuSy gene family may have different
2 functional roles and in sugarcane this was observed as different expression levels related to different TFBS
328 identified. We identified five different top-ranked TFBS (with the highest score) in the ScSuSy1-5 members.
4
329 Three of them are related to auxin and abscisic-acid hormone signaling (ScSuSy1, 3, 5). For ScSuSy1 genes,
7
330 the TFBS analysis predicted the motif wATATATATw (MA1184.1) that is associated with RVE1, a morning-
9
331 phased transcription factor integrating the circadian clock and auxin pathway genes that bind to the evening
11
332 element (EE) of promoters [55]. For ScSuSy2 genes, we found the motif GACrAATryA (MA1374.1) that is
13
333 associated with IDD which regulates photoperiodic flowering by modulating sugar transport and metabolism
15
334 [56]. For ScSuSy3 genes, we found the AyACTAGTrT (MA0930.1) motif in 64% of its SP80-3280 copies and
18
335 in all copies in the *S. spontaneum* and R570 monoploid genomes. It is associated with ABA-responsive
20
336 elements (ABRE) that regulate stress response via ABA signaling. For ScSuSy4 genes, we found the
22
337 TAGyAynTTT (MA1012.1) motif that is probably involved in regulation of the photoperiod and vernalization
24
338 pathways. Finally, for ScSuSy5 genes, we found a CTGCTAGCAG (MA0564.1) conserved motif exclusively
27
339 for ScSuSy5 genes in SP80-3280. This motif allows binding with an element associated with ABI3, which
29
340 participates in abscisic acid (ABA)-regulated gene expression. Previous studies from our group had already
31
341 pointed out ABA- and sucrose-induced genes associated with higher sucrose content in sugarcane [57].
33
342

343 On the other hand, for ScPAL I genes, the TFBS analysis predicted an ArCAyATnTG (MA0930.1)
38
344 element, which is associated with ABF3, a transcription factor involved in ABA and stress responses and
40
345 acting as a positive component of glucose signal transduction. For ScPAL III genes, we found the element
42
346 GGTCsGGcKc (MA0992.1), an element associated with AP2/ERF, a transcription factor involved in the
44
347 regulation of gene expression by stress factors and by components of stress signal transduction pathways. For
46
348 ScPAL Va genes, we found the element TCTAAAGTTT (MA0064.1), which is associated with PBF, a
48
349 transcription factor involved in ABA, stress response and components of stress signal transduction pathways.
50
350 Finally, for ScPAL Vb genes, we found the motif GCCGGAACGG (MA1009.1). This element is associated
52
351 with ARF3, a transcription factor involved in auxin and ABA-regulated gene expression. In summary, our
54
352 results corroborates reported findings [57] which reveal that PAL genes were induced by ABA.
56
58
59
60
61
62
63
64
65

353 Fewer TEs were identified in SP80-3280 in comparison to the two other published genome [13,14] and
354 few cultivar specific amplifications were observed as all previously described TE families were identified.

355 Around 81% of the SNVs identified were bi-allelic with only 0,97% being tetra-allelic. The predominance
356 of biallelic variations supports the theory that sugarcane experienced allotetraploidy shortly after divergence
357 with sorghum ca. 3.8~4.6 MYA [58], creating two ‘subgenomes’. Further, autotetraploidization after
358 *Saccharum* speciation ca. 3.1~3.8 MYA may have contributed to allelic richness within each sugarcane
359 ‘subgenome’, with occasional exchanges between subgenomes likely. Recent published results from Vieira et
360 al. [59], demonstrate that sugarcane meiotic chromosomes behave as bivalents, supporting these inferences.

361 In an attempt to organize the contigs, we allocate them in 6 clusters using MCA with hierarchical
362 clustering of the sequences. MCA results suggest that 4 out of 6 clusters correspond to single chromosomes in
363 *S. spontaneum* and R570. On the other hand, clusters 3 and 4, which contain multiple chromosomes, include
364 those in which chromosomal rearrangement events were demonstrated in comparison to sorghum: SsChr5,
365 SsChr6 and SsChr7 from *S. spontaneum* [14] and six R570 hom(oe)ology groups HG5-HG10 [13].

368 CONCLUSION

369 The singular challenge associated with sugarcane breeding makes genomic tools and genome assembly of
370 a polyploid interspecific hybrid of especially high value. Its large autopolyploid genome, predominantly clonal
371 propagation, and need for extensive phenotyping to determine breeding values, have contributed to the
372 relatively slow (~1% per year at most) rate of progress in improvement of sugarcane [60] and perhaps other
373 autopolyploids. The demonstration that most of its many homo(eo)logs are expressed, often with tissue-
374 specificity, and that transcription factor binding sites and TE insertions differ among homo(eo)logs, suggests
375 complex constraints that may necessitate unusual richness of information to make effective decisions about
376 selecting some homo(eo)logs alleles at the expense of others in autopolyploid breeding populations. These
377 principles may apply widely to many plants with large polyploid genomes that include many of those most
378 efficient at converting solar radiation to biomass.

379 The present work discloses a large collection of gene-space homo(eo)logs diversity, taking advantage of
380 novel sequencing technologies, adding over 3Gb of sequence not previously reported, in addition to genome

381 annotation, data mined homo(eo)logs, and explored regulatory regions. The resolved gene-space of the
382 sugarcane genome is a fundamental step towards a high-quality chromosome resolved assembly from a current
2 commercial hybrid. The genome sequence released for this interspecific polyploid supports its recent
383 allotetraploid nature, reveals differences in promoter regions associated to differentially expressed genes and
4
384 transposable elements contributing to fine tuning of the rich diversity in a genome that is otherwise highly
6
385 syntenic with its close relative, sorghum.
7
386

11
1387

13
1388

15
1389

METHODS

17
18
1390

19
20
21
22
2391

Plant material

24
25
26
27
28
29
30
31
32
3392

Leaves from SP80-3280 were collected and frozen in liquid nitrogen. Genomic DNA was extracted using
24 DNeasy Plant Mini Kit (Qiagen) following the standard protocol. DNA integrity was analyzed using the
25 Agilent High Sensitivity DNA Analysis Kit (Agilent Technologies) and Agilent 2100 Bioanalyzer Instrument.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

395 Quantification was done using Quant-it™ PicoGreen® dsDNA Assay Kit (ThermoFisher Scientific) and
396 SpectraMax M2 microplate reader (Molecular Devices).

36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Sequencing Illumina Long-reads and Assembly

We used Illumina Synthetic Long-read sequencing technology, which provides very accurate long reads with
40 a mean read length of roughly 5 kb, thus being able to represent polymorphisms across all copies of
41 chromosomes. Genomic DNA was sheared into 5-10 kb fragments and diluted in 384-well plates. DNA
42 fragments were ligated with PCR primers and specific sequences, which identify the 5' and 3' ends. The
43 fragments were amplified, fragmented and barcoded to create 26 TruSeq Synthetic Long-Read DNA libraries.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

404 The short fragments created in the second step of fragmentation were pooled and sequenced on the HiSeq
405 instrument at the Illumina Service Genome Network. The reads from each of the 384 wells were pre-processed
406 to correct sequencing and PCR errors. Contigs were produced from the paired-end information and further
407 scaffolded together to resolve repeats and fill in gaps. More details on the informatics pipeline for short read
408 scaffolding into long reads are available in the Fast Track Services Long Reads Pipeline User Guide [61].

409 To assemble sequences we used a two step approach: *i*) the Celera Assembler [62] (CA) was used for overlap
410 computation and layout building; *ii*) the *tig-sense* module of the HBAR-DTK (Hierarchical-Based AssembleR
2
411 Development ToolKit) from Pacific Biosciences [63] was used to construct consensus sequences. This was
4
412 motivated by the fact that the CA, which uses the overlap-layout-consensus method, is more robust than *de*
5
6
7
413 *Bruijn* graph approaches. However, some adjustments needed to be made. CA, designed for Sanger reads, only
9
414 accepts quality scores between 0 and 40. Since synthetic long reads are very accurate and some of the base
11
415 qualities exceeded this upper bound, we modified the quality scores of our long-read data to allow them to be
13
416 appropriately parsed. The consensus module was also adapted for the analysis of big complex genomes. The
14
15
16
417 substantial number of contigs generated initially (roughly 450.000, half of them singletons) incurred in several
17
18
418 files in a folder that hindered I/O operations. So, we *i*) modified *tig-sense* to automatically create subdirectories
20
21
419 that contained not more than a thousand contig FastA files, reducing delays for file lookup; *ii*) divided contig
22
23
420 processing into non-singletons and singletons, prioritizing non-singleton contigs; and *iii*) created a work
24
25
421 history so that the program could be resumed after a halt. Overall, these modifications allowed us to reduce
26
27
422 the running time of the consensus pipeline by one or two orders of magnitude.
28
29

30 31 32 33 **Sequencing BAC clones and assembly**

34
425 A total of 780 independent BACs were sequenced using Roche454 and PACBio sequencing technologies.
35
36
426 Each BAC clone was tagged with a unique barcode and sets of 12 BACs were pooled in one gasket. We
37
38
427 assembled BACs individually as described [64] and obtained a total of 49.6 Mbp of assembled sequence, with
39
40
41
428 a mean length of 107 kbp. The BAC data includes 317 R570 BACs¹⁸, 116 additional R570 BACs and 347 from
42
43
429 SP80-3280.
44
45

46 47 48 49 **Assembly Validation**

50 51 ***Comparison with Sugarcane BACs***

52
433 Assembled contigs were aligned against a set of 780 BACs with BWA mem. Alignment data was processed
53
54
55
434 for coverage with the aid of Samtools (v1.1) and Bedtools (v2.25) and selected matches were at least 10 kbp
56
57
435 long and covered 90% or more of the contig. Additionally, the unassembled synthetic long reads were aligned
58
59
60
61
62
63
64
65

436 to the same set of BACs, to check for discrepancies among contigs and long reads, which could be indicative
437 of regions that were not assembled.

2
438

4
439

Comparison with Sorghum CDS

6
7
440

The set of 39,207 annotated sorghum coding sequences (CDS), release version v2.1, were downloaded from
441 Phytozome [65]. These were aligned against the assembled contigs with BLASTn (v2.2.30+) using default
442 parameters. For each sorghum CDS, we identified the longest fraction of the coding sequence contained within
443 a single unitig. Only hits with at least 80% identity at the nucleotide level were considered for computing
444 coverage.

16
17
18
19
445

20
21
2446

Comparison with CEGMA

22
23
2447

A total of 248 Ultra-conservative core eukaryotic genes classified by Korf Lab [22] were assessed in
448 our sugarcane assembly with ‘-g’ and other default options of CEGMA v2.5.

25
26
27
2449

29
30
3450

Comparison with BUSCO

31
32
3451

Assembly completeness was assessed by searching for the 1,440 core genes from the Plantae lineage of
452 Benchmarking Universal Single-Copy Orthologs (BUSCO) [23]. BUSCO performs gene prediction and
453 orthogonality assessment using Augustus [66] and HMMER3 [67]. Since these steps demand huge resources,
454 we parted sugarcane contigs (4.3Gbp) into six groups with similar volume and processed BUSCO in parallel.
455 After we merged results, we applied orthogonality assessment algorithm once again as thresholds that BUSCO
456 exploits to discern actual single copy orthologs from paralogs.

45
46
4457

47
48
4458

Comparison of the mitochondrial and chloroplast genomes

49
50
4459

To reconstruct the SP80-3280 mitochondrial and chloroplast genomes, we have used as reference the complete
460 genomes of *Saccharum* hybrid chloroplast (NC_005878.2) [68] and the *Saccharum officinarum* mitochondrial
461 chromosome 1 (LC107874.1) and chromosome 2 (LC107875.1) [25], downloaded from NCBI. The SP80-
462 3280 genome contigs were aligned using BLASTN against their respective references and the best hits were
463 selected based on cutoff E-value $\leq 1.0E-15$, with contig coverage $\geq 90\%$ and identity $\geq 70\%$. The BLASTN

59
60
61
62
63
64
65

464 alignment results identified 2,482 and 909 contigs for the two mitochondrial chromosomes, respectively; and
465 51,768 contigs for the chloroplast genome. To reconstruct the consensus sequences and do the genome
2
466 annotation we have used the CLC Genomics Workbench tools [69]. Using the CLC Tools and the Genome
4
467 Finishing Module, the selected contigs were aligned to their respective references and consensus sequences
6
468 extracted, filling the gaps with N's. The reconstructed consensus sequence aligned against the chloroplast
7
9 genome presented 100% of coverage and identity. Alignment against mitochondrial chromosomes 1 and 2
10
1470 presented over 99% of coverage and identity. The consensus sequences were annotated using their respective
13
1471 NCBI references with the CLC tool "Annotate from Reference", where all genes, tRNAs rRNAs and
15
1672 miscellaneous features were totally transferred.
17

18
1973

20

2174

Genome Annotation

22

2375

Gene prediction

24

2576

Contigs were annotated using a pipeline developed in house, previously used for BAC annotation.

26

2777

Transposable element (TE) discovery and masking was done using LTR harvest, LTR digest, CrossMatch

28

3078

against *Utricularia gibba* TE DB and RepeatMasking [70] of Viridiplantae [71] and previously known

31

3279

sugarcane TEs [45].

33

3480

Genes were discovered and annotated using masked contig sequences. *De novo* predictions were done with

35

3681

Augustus [66], Glimmer HMM [72], GeneMark HMM [73], SNAP and PASA [74] with rice models and

37

3882

sugarcane EST and RNAseq data. Alignments were also generated against reference protein DBs (sorghum,

39

4083

known sugarcane and Phytozome) using Exonerate [75] and BLAST [76] (v2.2.30+). Both *de novo* and

41

4284

alignment evidence were used for consensus annotation with EVIDENCEModeler [77] with greater weight given

43

4485

to experimental and alignment information. Functional assignment was derived from protein DB best hits and

45

4686

InterProScan 5 [78] results.

47

4887

49

50

5188

GeneOntology annotation

52

5389

For functional annotation of predicted proteins from SP80-3280, all sequences were aligned to UniRef50

54

5590

clusters, a dataset of representative sequences clustering high similarity proteins from UniProtKB [30], using

56

5791

BLASTP (v2.2.30+, *-evalue 1e-5*). Sequences that fail to align in this first approach were also searched against

58

59

60

61

62

63

64

65

492 the RefSeq non-redundant protein database. Gene Ontology mapping and annotation of sequences with
493 positive BLAST results was performed using Blast2Go framework [79].

2
494
4

495 **Reference-guided RNAseq Assembly**

496 We used Trinity version 2.0.6 for reassembly of the Sugarcane ORFeome [29] using the genome as a reference,
497 with a minimum contig length of 250 bp (genome_guided_max_intron 3000, genome_guided_min_coverage
498 5, genome_guided_min_reads_per_partition 10) to identify transcript models. SP80-3280 RNA-seq reads from
499 3 tissues (leaves and immature and intermediate internodes) were used for alignment against the reference
500 genome and partitioned into read clusters, which were then individually assembled using Trinity genome-
501 guided methods. Trinity and genome-guided methods used a fixed k-mer size of 25nt. In this new assembly,
502 269,050 genes and 275,807 transcripts were recovered. The quantity of transcripts recovered by the reference
503 guided-assembly was higher, and thus closer to the number of predicted genes (374,774), than the *de novo*
504 assembly.

25
504
26
27
505

506 **Identification of Gene Copies and Count Estimation**

507 We downloaded the *Sorghum bicolor* genome assembly v2.1 from Phytozome and took 2,051 single copy
508 genes according to Han *et al.* [80], which were also present as single copies in the genomes of *Oryza sativa*
509 and *Brachypodium distachyon*. We aligned the coding sequences of these sorghum genes to the coding
510 sequences of predicted sugarcane genes from the SP80-3280 assembly, using the BLASTn (v2.2.30+, -*evaluate*
511 *le-6*). We filtered alignments with at least 80% nucleotide identity, covering at least 70% of both the sugarcane
512 and sorghum sequences. For sorghum coding sequences with multiple hits to sugarcane genes, we further
513 required that each hit had at least 80% identity, based on Wang *et al.* [46]. Sugarcane gene models aligned to
514 the same single copy sorghum gene were denoted as putative homo(eo)logues. Finally, we counted the number
515 of copies for each gene.

516 We clustered all gene copies based on each single copy sorghum gene to get estimates of sequence
517 differentiation. We aligned the coding sequences for each pairwise combination in each gene cluster, using
518 BLAT v35 [81] (*-minIdentity=0 -minScore=60*), disregarding clusters with more than 16 putative
519 homo(eo)logs. Next, we parsed the alignments to obtain estimates of copy differentiation considering both

59
60
61
62
63
64
65

520 SNPs and INDELS. We gathered distance estimates from all pairs, from all clusters, to obtain dissimilarity
521 distributions.

2
522

4

523 **Gene copy characterization**

6

524 *Upstream region analysis*

7

525 We also assessed the dissimilarity levels of regions upstream (potential promoter regions) of the predicted
526 sugarcane gene copies. We initially collected three different sequence ranges (100 bp, 500 bp and 1000 bp)
527 upstream of the predicted gene start site. Next, we aligned these upstream sequences for each pairwise
528 combination in each cluster, again using BLAT v35 [81] (*-minIdentity=0 -minScore=30*). Finally, for each
529 distance range, we parsed the alignments and computed the dissimilarity level considering both mismatches
530 and gaps. To avoid partial alignments of the upstream sequences, only alignments up to 20% shorter or longer
531 than the expected sequence length were considered.

24
25
532

26

533 *Insertions and Deletions between gene copy Coding Sequences*

27

534 To investigate the occurrence of frameshift mutations between gene copies, we built multiple alignments of
535 the coding sequences of putative copies, for each cluster, with MUSCLE v3.8.31 [83], using default
536 parameters. We then computed the length distribution of insertions and deletions in the coding sequences, to
537 differentiate between frame-preserving and frameshift indels. We parsed the CDS alignment for each pairwise
538 combination of gene copies and counted the number of occurrences of gaps of a given length. We then pooled
539 counts from all copy combinations to get a joint estimated distribution.

42
43
540

44

541 *Tissue-Specific Homo(eo)logs Expression Analysis*

45

542 We used RNA-Seq data [29] from leaves (*L*), immature (*II*) and intermediate (*I5*) internodes of SP80-3280 to
543 find the expression of putative tissue-specific gene copies. These reads were initially aligned to the Sugarcane
544 genome assembly using TopHat2 [84] version 2.0.9 (*library-type fr-firststrand*). We allowed reads to be
545 aligned to up to 20 contigs of the genome assembly to identify alignments to different homo(eo)logs (*--max-*
546 *multihits 20*) and supplied TopHat2 with the putative homo(eo)logs' annotation as a GTF file (*--GTF*
547 *CDSMapping-homo(eo)logs.gtf*), in order to direct TopHat2 to align the reads to this transcriptome first.

59
60

61

62

63

64

65

548 Besides the *TopHat2* alignment, we used the RSEM tool *rsem-calculate-expression* (version 1.2.31) to quantify
549 the expression of predicted genes (*bowtie2*, *fragment-length-mean*, *fragment-length-sd* and *calc-ci*
2 parameters). An in house perl script was used to estimate the mean length and standard deviation for each
550 RNA-seq library. The main output of *TopHat2* BAM formatted file [85] *accepted_hits.bam* was used with
4
551 *RSEM* to estimate the transcriptome expression profile. We developed in-house perl and R language (version
7
552 3.3.2) scripts to find the number of putative expressed homo(eo)logs for each single copy Sorghum gene, using
9
553 the information from *genome annotation* file (GFF format), showing the gene structure, the transcriptome
11
554 annotation and respective TPM (Transcript Per Million) abundance. The previous information allowed the
13
555 creation of the homo(eo)logs GFF file. We also applied TopHat2 to find the number of putative homo(eo)logs
15
556 expressed only in *antisense* orientation, using the same protocol described above, and the *antisense* reads of
17
557 RNA-Seq previously identified by Nishiyama *et al.* [29].
18
20
21
22

560 **ScSuSy and ScPAL gene family analysis**

561 We used the sugarcane and sorghum SuSy protein sequences reported by Zhang *et al.* [35] as query for a
27
28
29
30
31
32
33
34
562 tBLASTn (v2.2.30+) search in the predicted proteins from SP80-3280, *S. spontaneum* [46] and R570 genome
35
36
37
563 assemblies [13]. Putative SuSy genes were then filtered by query coverage $\geq 80\%$ of at least one of the five
38
39
564 ScSuSy from Zhang *et al.* [35] and by PFAM [86] domain search, considering only those containing both the
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
conserved sucrose synthase and glucosyl-transferase 1 domains.

66 Based in BLAST and keyword search in two databases (Plant GDB, <http://www.plantgdb.org/> and Phytozome
67 [65]) we found 8 different PAL genes in the sorghum genome, the same number previously reported [87]. For
68 sugarcane, PAL genes were retrieved from an EST Cell Wall catalogue [41], which was used as query together
69 with sorghum PAL genes for a BLASTn (v.2.2.30+) search to identify PAL genes in the predicted proteins
70 from *S. spontaneum* [47] and R570 genome assemblies [13]. Putative PAL genes were then filtered by query
71 coverage $\geq 80\%$ of the sorghum PAL genes and by PFAM [86] domain search, considering only those
72 containing the Aromatic amino acid lyase domain. Also, sequences not containing the PAL conserved amino
73 acid motif Ala-Ser-Gly [88,89] and an essential Tyr110 [90] were excluded.

74 For both SuSy and PAL, nucleotide sequences (CDS) were aligned with *clustalw* [91] software in MEGA 7.0
75 [92] and maximum likelihood trees were constructed with default parameters except for 1,000 bootstraps and

576 Gaps/missing data treatment “*use all sites*”. Expression heatmap was constructed using log₂ transcript per
577 million (TPM) from previous RNA-seq data [29].

2
578
4

579 **Cell wall-related genes**

580 For the identification of cell wall-related genes in the sugarcane genome we used the Sugarcane SAS Cell Wall
581 catalogue [41] as a reference. The search was carried out using tBLASTn (v2.2.30+, *-evaluate 1e-6*). These were
582 manually re-annotated to produce a sugarcane cell wall catalogue with 3,054 sequences, classified in 10 cell
583 wall categories.

14
15
16
584
17

585 **Transcription Factor analysis**

586 For the identification and classification of sugarcane predicted proteins into transcription factor (TF) families,
587 we used the classification rules and tools described in GRASSIUS [44]. The search was carried out using
588 HMMER v3.1b1 [93] and all significant HMM hits with *e*-value smaller than $1e^{-3}$ were kept.

25
26
27
28
289
29

590 **Promoter region analysis**

591 *Transcription Start Site (TSS) and promoter region classification*

592 We evaluated promoter regions of genes associated with cell wall and sugar metabolism, ScPAL and ScSuSy,
593 respectively, as described above. A total of 47 ScPAL and 44 ScSuSy was used. To extract the candidate
594 promoter region, we selected, when available, up to 1500 nt upstream from the annotated start position of the
595 gene, consisting of a core promoter (500 nt upstream of the start position) and proximal promoter (1000 nt
596 upstream of the core promoter). Next, we used TSSPlant [94] to predict the TSS of the genes and the type of
597 promoter (TATA-box, TATA-less). The software was set to report high score, sense only TSSs.

598 *Transcription Factor Binding Site (TFBS) in silico characterization*

599 The annotation of TFBSs in the proximal promoter regions was performed in two steps: *de novo* prediction of
600 TFBS motifs in smaller subsets of sequences and mapping the predicted TFBSs in the remaining promoter
601 sequences. Sequences were partitioned in 10 subsets: five ScPAL groups and five ScSuSy groups. We then
602 applied MEME [95] and MotifSampler [96], with default parameters, to each of these datasets to determine
603 putative TFBS motifs. Both were restricted to search for at most 6 motifs with 10nt or less. MEME candidates

59
60
61
62
63
64
65

604 were a subset of MotifSampler's. MotifSampler ran for 100 cycles; following the manual we selected, from
605 the 10 top-ranked motifs, the first 5 that occurred at least 10 times in the different cycles. Each of the resulting
2
606 35 candidate motifs was searched in the JASPAR public database [97], with partial positive matches for all of
4
607 them.

608 To evaluate the significance of the motifs we measured their frequency in promoter regions of each of the
7
9
1609 original gene families and compared them with the frequency of each of these motifs in the promoter regions
11
1610 of the other SP80-3280 predicted genes. We also mapped the motifs of each ScSuSy and ScPAL gene family
13
1611 respectively in the promoter region of the ScSuSy and ScPAL genes from *S. spontaneum* and R570. Candidate
15
1612 motifs were mapped with MotifLocator [96]. For characterizing background sequences, we trained a first order
17
18
1613 Markov chain [96] trained on SP80-3280 coding regions that were previously shuffled using the fasta-shuffle-
20
21
1614 letters tool [95]. The parameters were set to full match of the motif in the target sequence and score 95% above
22
23
1615 of the background.

2617 **Co-expression analysis**

2618 A field experiment was conducted at the Agricultural Sciences Center of the Federal University of São Carlos
31
32
2619 in Araras (22°21'25''S and 47°23'3''W) in the state of Sao Paulo, Brazil. Trial plots of SP-3280 consisted of
33
34
2620 four rows of 10 m long and spaced 1.35m apart. The field experiment was initiated in October 2012 and
36
37
2621 extended up until November 2013, representing the conditions under which “one-year” sugarcane crops are
38
39
2622 cultivated. Aiming to carry out observations throughout growth and development, tissue samples of the +1
40
41
2623 leaves (L1) and upper (I1), immature (I5) and mature (I9) internodes were collected from two plots (two
42
43
2624 technical replicates) after 4, 8, 11 and 13 months of planting.

45
2625 RNA was extracted for four biological replicates, two from each plot, using the TriZol method, treated with
47
48
2626 DNase I and purified. A pool of samples from leaves and a pool of internodes was used as a 'reference sample'
49
50
2627 for hybridization experiments on a customized 4 × 44 K oligoarray (Agilent Technologies) for sugarcane
51
52
2628 (CaneRegNet), conducted following the recommendations proposed by Lembke et al. [97]. The oligoarrays
53
54
2629 were read using the GenePix 4000B scanner device (Molecular Devices) and the fluorescence data was
56
57
2630 processed by Feature Extraction software 9.5.3 (Agilent Technologies).

631 Log2 transformed expression data was used for discovery and the analysis of co-expression modules,
632 on CEMiTool R package [98]. The adjacency matrix was calculated by estimating the Spearman's correlation
2
633 coefficient between all pair of genes and raised to a soft thresholding power (β) of 14. TopGO R package [99]
4
634 was used for gene ontology enrichment analysis for each module and node and edge files were generated for
6
635 use with the Cytoscape network visualization program [100].
7
8

9 10 11 12 **SNP variants (SNVs) analysis compared to genic regions in *Sorghum bicolor***

13
14
15 The 450,609 sugarcane contigs (183,322 singletons and 267,287 unitigs) were aligned to the sorghum genome
16
17 sequence [48] using the BWA MEM v0.7.10 [101] and reads with alignment score larger than 20 were used
18
19 for variant calling. SNVs were called using samtools v1.1 and bcftools v1.1 [85]. Using in-house python
20
21 scripts, extracted SNVs were screened when sugarcane contigs were located on the genic regions of the
22
23 sorghum genome and two or more sugarcane contigs were aligned to the same sorghum gene. Then, the number
24
25 of SNVs in each gene was counted according to four-base changes.
26
27 SNVs that are homozygous in sugarcane were extracted for further analysis. Large-effect SNVs were identified
28
29 as those mapped to coding regions, splicing sites, stop codons and transcription initiation sites.
30
31

32 33 34 **Functional Enrichment Test**

35
36
37 *Arabidopsis* GO-slim gene annotation was used for functional enrichment analysis. GO-slim terms were
38
39 assigned to sugarcane genes based on sequence similarity inferred from best BLASTp (v2.2.30+) hit. We used
40
41 a binomial distribution based on the proportion of a GO-slim term among all annotated genes in the sorghum
42
43 genome as the null distribution. The binomial test was used to assess functional enrichment, with a significance
44
45 threshold of $p > 0.05$.
46
47

48 49 50 **Conserved Synteny Blocks**

51
52 DNA sequences for all CDSs from *S. spontaneum* [47], R570 [13], *Sorghum bicolor* [102] and SP80-
53
54 3280 were aligned using the BLASTn program. Results from BLAST searches, with e-value $\leq 10^{-5}$, were
55
56 parsed using an in-house Python script to filter alignments covering at least 70% of the length of both the query
57
58 and hit sequences. A second filter, requiring at least 80% identity was also applied and the resulting pairs of
59
60

659 queries and hit sequences were classified into putative orthologous groups using the union-find algorithm. We
660 selected putative orthologous groups present in all three organisms but with only one *Sorghum* gene to be used
2
661 as markers to detect blocks of conserved gene order (syntenic blocks) in comparisons of SP80-3280 and *S.*
4
662 *spontaneum* against the genome of *S. bicolor*, thus avoiding the complications of a direct comparison of the
6
663 two polyploid genomes (**Additional file 1: Fig. S8**). Another Python script was used to detect the syntenic
8
664 blocks in both *Saccharum* genomes and to count the number of syntenic blocks in each contig.
10

665 666 **Chromosome Synteny Multiple Correspondence Analysis with Clustering**

667 We performed a multiple correspondence analysis (MCA) with clustering of the best local alignment hit of
17
668 masked contigs. Input data were the 450,609 contigs of the sugarcane synthetic long read assembly and the
19
669 masked genomic sequences of *S. spontaneum* [47] and R570 [13]. We used the masked sugarcane contig
22
670 sequence produced by the annotation pipeline, excluding 69,879 sequences that were fully masked.
24

671 The contigs were aligned to the grass genomes using BWA-SW v0.7.12-r1044 [101]. We used an in-house
26
672 Perl 5 script to retrieve the highest scoring hit for each contig and generate a table for input into R v3.2.1 [82].
28
673 This table contained the chromosome hit, if any, for each contig against each reference genome.
31

674 We then used the FactoMineR R package v1.31.3 [103], along with the missMDA missing data handling
33
675 auxiliary package v1.8.2 [104]. We performed MCA with these data, *i.e.*, chromosome hit number information
36
676 for each contig was treated as a set of categorical variables and represented in the two principal component
38
677 dimensions. This was followed by hierarchical clustering in these two dimensions, as well as figure rendering,
40
678 using the Hierarchical Clustering on Principal Components (HCPC) function of FactoMineR.
42

679 In order to identify the correspondence between *S. spontaneum* and R570 chromosomes and SP80-3280
44
680 clusters, protein sequence alignment between the cultivar variety and the ancestor and R570 was performed
47
681 with BlastP considering an e-value threshold of $1e^{-5}$. The best hit with a minimum query coverage of 90%
49
682 was selected for visual representation of the alignment results with Circos plot.
51

683 684 685 **ADDITIONAL FILES**

686 **Additional file 1.doc contains Supplemental Figures S1 to S10**

687 **Additional file 2.xls contains Supplemental Tables S1 to S13**

688

2

689 **DECLARATIONS**

4

690

6

691 **List of abbreviations**

8

9

692

10

11

693 CEGMA: Core Eukaryotic Genes Mapping Approach

13

694 BUSCO: Benchmarking Universal Single-Copy Orthologs

15

695 ESTs: expressed sequence tags

17

18

696 CDS: coding sequences

19

20

697 SuSy: Sucrose Synthases

21

22

698 PAL: Phenylalanine ammonia-lyase

24

25

699 CCR: Cinnamoyl-CoA reductase

26

27

700 COMT: Caffeic acid 3-O-methyltransferase

28

29

701 4CL: 4-coumarate-CoA ligase

31

32

702 TFBSs: Transcription Factor Binding Sites

33

34

703 TE: transposable elements

35

36

704 MCA: Multiple correspondence analysis

37

38

705 I2C-2: R gene locus

39

40

706 *dog1*: (delay of germination 1

42

43

707 ABRE: ABA-responsive elements

44

45

708 ABA: abscisic acid

46

47

709

48

49

710 **Consent for publication:** Not applicable

51

52

711

53

54

712 **Availability of data and material**

55

56

713 Genomic data is publicly available at NCBI under GenBank Bioproject PRJNA431722. Contig sequence and

57

58

714 annotation are also available in a genome browser framework at <http://sucest-fun.org/>. The microarray data

60

61

62

63

64

65

715 have been deposited in NCBI's Gene Expression Omnibus and are accessible through GEO Series accession
716 number GSE124990.

2
717
4

718 **Competing interests**

719 The authors declare that they have no competing interests.

9
720

11
721

721 **Funding**

13
722

722 This work was funded by State of São Paulo Foundation and Microsoft Research (FAPESP grant n°
15 2012/51062-3) and State of São Paulo Foundation (FAPESP grants n° 2014/50921-8, 2008/52146-0 and
16 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
17 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
18 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
19 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
20 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
21 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
22 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
23 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
24 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
25 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
26 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
27 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
28 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
29 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
30 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
31 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
32 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
33 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
34 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
35 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
36 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
37 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
38 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
39 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
40 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
41 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
42 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
43 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
44 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
45 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
46 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
47 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
48 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
49 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
50 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
51 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
52 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
53 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
54 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
55 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
56 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
57 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
58 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
59 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
60 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
61 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
62 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
63 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
64 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science
65 2008/52074-0) under the BIOEN Program. Additional funding included awards from the National Science

15

16
17

18
19

20
21

22
23

24
25

26
27

28
29

30
31

32
33

34
35

36
37

38
39

40
41

42
43

44
45

46
47

48
49

50 **Authors' contributions**

51

52 Project leaders: GMS, MAVS and DH;

53

54 Sample collection and DNA extraction: CGL;

55

56 Genome sequencing and assembly: HL, MCS, GRAM, RP and BD;

57

58 Genome assembly supervision: DH;

59

60 Genome annotation: MAVS, GJW, MYNJ and FTC;

61

62
63
64
65

744 Saccharum spontaneum genome assembly: JZ, XZ, QZ and RM;
745 BWA-SW analysis: GJW;
2
746 BAC sequencing and assembly: MAVS, GJW, GTR, HB and SV;
4
747 Synteny analysis: AMD, RFS and GGS;
6
748 Reference-guided RNAseq Assembly: MYNJ;
8
9
749 Tissue-Specific Allelic Expression Analysis: MYNJ, CGL and PMA;
10
11
750 Phylogeny analysis: SSF and ALD;
13
751 SP80-3280 growth and maturation experiment: MSC, GMS, CGL and ALD
15
752 Co-expression analysis: ALD
17
18
753 Regulatory region analysis (TE and TFBS): MAVS, MMO, AMD, GMS, CTH and ALD;
19
20
754 SNP variants (SNVs) analysis: CK, HG and AP;
22
755 Organization and management of the author's contributions: CGL, ALD, GMS and MAVS;
24
756 Data availability (NCBI and Sucest-fun): FTC;
26
757 All authors have read and approved the final version of the manuscript.
28
29

30

31

32 **Acknowledgements**

33

34 We are indebted to Andreia Prata, Vania Sedano, Nathalia de Setta, Joni Lima, Marcos Buckeridge, Eveline

35

36 Tavares, Katia Scortecci, Anete Pereira de Souza, Sonia Vautrin and Hélène Bergès for contributions in BAC

37

38 library construction, BAC selection or sequencing. We are indebted to the Sugarcane Genome Sequencing

39

40 Initiative for useful discussions.

42

43

44

45 **REFERENCES**

46

47

48 1. FAOSTAT. Production/Crops, Food and Agriculture Organization of the United Nations - Statistics Division
49 [Internet]. 2018. Available from: <http://www.fao.org/faostat/en/#home>

50

51 2. Long SP, Karp A, Buckeridge SC, Davis SC, Jaiswal D, Moore PH, et al. Feedstocks for biofuels and bioenergy.
52 Bioenergy Sustain Bridg Gaps [Internet]. Paris Cedex: Scientific Committee on Problems of the Environment
53 (SCOPE); 2015. p. 302–347. Available from: http://bioenfapesp.org/scopebioenergy/images/chapters/bioen-scope_chapter10.pdf

54

55 3. Kline KL, Msangi S, Dale VH, Woods J, Souza GM, Osseweijer P, et al. Reconciling food security and
56 bioenergy: priorities for action. GCB Bioenergy. 2017;9:557–76.

57

58 4. Goldemberg J. Ethanol for a Sustainable Energy Future. Science. 2007;315:808–10.

59

60

61

62

63

64

65

- 776 5. Jaiswal D, De Souza AP, Larsen S, LeBauer DS, Miguez FE, Sparovek G, et al. Brazilian sugarcane ethanol as
777 an expandable green alternative to crude oil use. *Nat Clim Change*. 2017;7:788–92.
- 1
778 6. Souza GM, Ballester MVR, de Brito Cruz CH, Chum H, Dale B, Dale VH, et al. The role of bioenergy in a
779 climate-changing world. *Environ Dev*. 2017;23:57–64.
- 4
780 7. Souza GM, Victoria RL, Joly CA, Verdade LM. *Bioenergy & sustainability: bridging the gaps*. Paris Cedex:
781 Scientific Committee on Problems of the Environment (SCOPE); 2015.
- 8
782 8. Souza GM, Filho RM. *Industrial Biotechnology and Biomass: What Next for Brazil’s Future Energy and
783 Chemicals?* *Ind Biotechnol*. 2016;12:24–5.
- 11
784 9. Vilela M de M, Del-Bem L-E, Van Sluys M-A, de Setta N, Kitajima JP, Cruz GMQ, et al. Analysis of three
785 sugarcane homo/homeologous regions suggests independent polyploidization events of *Saccharum
786 officinarum* and *Saccharum spontaneum*. *Genome Biol Evol*. 2017;evw293.
- 16
787 10. Jannoo N, Grivet L, Seguin M, Paulet F, Domaingue R, Rao PS, et al. Molecular investigation of the genetic
788 base of sugarcane cultivars. *Theor Appl Genet*. 1999;99:171–84.
- 19
789 11. D’Hont A. Unraveling the genome structure of polyploids using FISH and GISH; examples of sugarcane
790 and banana. *Cytogenet Genome Res*. 2005;109:27–33.
- 22
791 12. Thirugnanasambandam PP, Hoang NV, Henry RJ. The Challenge of Analyzing the Sugarcane Genome.
792 *Front Plant Sci* [Internet]. 2018 [cited 2018 Aug 23];9. Available from:
793 <http://journal.frontiersin.org/article/10.3389/fpls.2018.00616/full>
- 27
794 13. Garsmeur O, Droc G, Antonise R, Grimwood J, Potier B, Aitken K, et al. A mosaic monoploid reference
795 sequence for the highly complex genome of sugarcane. *Nat Commun* [Internet]. 2018 [cited 2018 Aug 16];9.
796 Available from: <http://www.nature.com/articles/s41467-018-05051-5>
- 32
797 14. Zhang J, Zhang X, Tang H, Zhang Q, Hua X, Ma X, et al. Allele-defined genome of the autopolyploid
798 sugarcane *Saccharum spontaneum* L. *Nat Genet*. 2018;50:1565–73.
- 35
799 15. Waclawovsky AJ, Sato PM, Lembke CG, Moore PH, Souza GM. Sugarcane for bioenergy production: an
800 assessment of yield and regulation of sucrose content. *Plant Biotechnol J*. 2010;8:263–76.
- 39
801 16. Goldemberg J, Coelho ST, Guardabassi P. The sustainability of ethanol production from sugarcane. *Energy
802 Policy*. 2008;36:2086–97.
- 42
803 17. Welbaum GE, Meinzer FC. Compartmentation of solutes and water in developing sugarcane stalk tissue.
804 *Plant Physiol*. 1990;93:1147–53.
- 45
805 18. Bonawitz ND, Chapple C. The genetics of lignin biosynthesis: connecting genotype to phenotype. *Annu
806 Rev Genet*. 2010/09/03. 2010;44:337–63.
- 49
807 19. Himmel ME, Ding SY, Johnson DK, Adney WS, Nimlos MR, Brady JW, et al. Biomass recalcitrance:
808 engineering plants and enzymes for biofuels production. *Science*. 2007/02/10. 2007;315:804–7.
- 52
809 20. Vettore AL. Analysis and Functional Annotation of an Expressed Sequence Tag Collection for Tropical Crop
810 Sugarcane. *Genome Res*. 2003;13:2725–35.
- 56
811 21. Riaño-Pachón DM, Mattiello L. Draft genome sequencing of the sugarcane hybrid SP80-3280.
812 *F1000Research*. 2017;6:861.

- 813 22. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic genomes.
814 Bioinformatics. 2007;23:1061–7.
- 1
815 23. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly
816 and annotation completeness with single-copy orthologs. Bioinformatics. 2015;31:3210–2.
- 4
817 24. Calsa Júnior T, Carraro DM, Benatti MR, Barbosa AC, Kitajima JP, Carrer H. Structural features and
818 transcript-editing analysis of sugarcane (*Saccharum officinarum* L.) chloroplast genome. Curr Genet.
819 2004;46:366–73.
- 9
820 25. Tsuruta S, Ebina M, Kobayashi M, Takahashi W. Complete Chloroplast Genomes of *Erianthus*
821 *arundinaceus* and *Miscanthus sinensis*: Comparative Genomics and Evolution of the *Saccharum* Complex.
822 Heinze B, editor. PLOS ONE. 2017;12:e0169992.
- 13
823 26. Nah G, Im J-H, Lim S-H, Kim K, Choi AY, Yook MJ, et al. Complete chloroplast genomes of two *Miscanthus*
824 species. Mitochondrial DNA Part A. 2016;27:4359–60.
- 17
825 27. Cardoso-Silva CB, Costa EA, Mancini MC, Balsalobre TWA, Canesin LEC, Pinto LR, et al. De Novo Assembly
826 and Transcriptome Analysis of Contrasting Sugarcane Varieties. Gibas C, editor. PLoS ONE. 2014;9:e88462.
- 20
827 28. Vicentini R, Bottcher A, Brito M dos S, dos Santos AB, Creste S, Landell MG de A, et al. Large-Scale
828 Transcriptome Analysis of Two Sugarcane Genotypes Contrasting for Lignin Content. Amancio S, editor. PLOS
829 ONE. 2015;10:e0134909.
- 25
830 29. Nishiyama MY, Ferreira SS, Tang P-Z, Becker S, Pörtner-Taliana A, Souza GM. Full-Length Enriched cDNA
831 Libraries and ORFeome Analysis of Sugarcane Hybrid and Ancestor Genotypes. PLOS ONE. 2014;9:e107351.
- 28
832 30. Suzek BE, Wang Y, Huang H, McGarvey PB, Wu CH, the UniProt Consortium. UniRef clusters: a
833 comprehensive and scalable alternative for improving sequence similarity searches. Bioinformatics.
834 2015;31:926–32.
- 33
835 31. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology: tool for the
836 unification of biology. Nat Genet. 2000;25:25–9.
- 36
837 32. Veeckman E, Ruttink T, Vandepoele K. Are We There Yet? Reliably Estimating the Completeness of Plant
838 Genome Sequences. Plant Cell. 2016;28:1759–68.
- 40
839 33. Nelson JC, Wang S, Wu Y, Li X, Antony G, White FF, et al. Single-nucleotide polymorphism discovery by
840 high-throughput sequencing in sorghum. BMC Genomics [Internet]. 2011 [cited 2018 Jan 26];12. Available
841 from: <http://bmcgenomics.biomedcentral.com/articles/10.1186/1471-2164-12-352>
- 44
842 34. Coleman HD, Yan J, Mansfield SD. Sucrose synthase affects carbon partitioning to increase cellulose
843 production and altered cell wall ultrastructure. Proc Natl Acad Sci. 2009;106:13118–23.
- 48
844 35. Zhang J, Arro J, Chen Y, Ming R. Haplotype analysis of sucrose synthase gene family in three
845 *Saccharum* species. BMC Genomics. 2013;14:314.
- 51
846 36. Persia D, Cai G, Del Casino C, Faleri C, Willemse MT, Cresti M. Sucrose synthase is associated with the cell
847 wall of tobacco pollen tubes. Plant Physiol. 2008;147:1603–18.
- 55
848 37. Brill E, van Thournout M, White RG, Llewellyn D, Campbell PM, Engelen S, et al. A Novel Isoform of Sucrose
849 Synthase Is Targeted to the Cell Wall during Secondary Cell Wall Synthesis in Cotton Fiber. Plant Physiol.
850 2011;157:40–54.

- 851 38. Sewalt V, Ni W, Blount JW, Jung HG, Masoud SA, Howles PA, et al. Reduced Lignin Content and Altered
852 Lignin Composition in Transgenic Tobacco Down-Regulated in Expression of L-Phenylalanine Ammonia-Lyase
853 or Cinnamate 4-Hydroxylase. *Plant Physiol.* 1997;115:41–50.
2
- 854 39. Rohde A. Molecular Phenotyping of the *pal1* and *pal2* Mutants of *Arabidopsis thaliana* Reveals Far-
855 Reaching Consequences on Phenylpropanoid, Amino Acid, and Carbohydrate Metabolism. *PLANT CELL*
856 *ONLINE.* 2004;16:2749–71.
4
5
6
- 857 40. Vanholme R, Storme V, Vanholme B, Sundin L, Christensen JH, Goeminne G, et al. A Systems Biology View
858 of Responses to Lignin Biosynthesis Perturbations in *Arabidopsis*. *Plant Cell.* 2012;24:3506–29.
7
- 859 41. Ferreira SS, Hotta CT, Poelking VG de C, Leite DCC, Buckeridge MS, Loureiro ME, et al. Co-expression
860 network analysis reveals transcription factors associated to cell wall biosynthesis in sugarcane. *Plant Mol Biol.*
861 2016;91:15–35.
10
11
12
13
14
- 862 42. Cunha CP, Roberto GG, Vicentini R, Lembke CG, Souza GM, Ribeiro RV, et al. Ethylene-induced
863 transcriptional and hormonal responses at the onset of sugarcane ripening. *Sci Rep [Internet].* 2017 [cited
864 2018 Aug 16];7. Available from: <http://www.nature.com/articles/srep43364>
15
16
17
18
19
- 865 43. Xu Z, Zhang D, Hu J, Zhou X, Ye X, Reichel KL, et al. Comparative genome analysis of lignin biosynthesis
866 gene families across the plant kingdom. *BMC Bioinformatics.* 2009;10:S3.
20
21
22
- 867 44. Yilmaz A, Nishiyama MY, Fuentes BG, Souza GM, Janies D, Gray J, et al. GRASSIUS: A Platform for
868 Comparative Regulatory Genomics across the Grasses. *PLANT Physiol.* 2009;149:171–80.
23
24
25
26
- 869 45. Domingues DS, Cruz GM, Metcalfe CJ, Nogueira FT, Vicentini R, de S Alves C, et al. Analysis of plant LTR-
870 retrotransposons at the fine-scale family level reveals individual molecular patterns. *BMC Genomics.*
871 2012;13:137.
27
28
29
30
- 872 46. Wang J, Roe B, Macmil S, Yu Q, Murray JE, Tang H, et al. Microcollinearity between autopolyploid
873 sugarcane and diploid sorghum genomes. *BMC Genomics.* 2010;11:261.
31
32
33
34
- 874 47. Zhang et al. Allele-defined genome of the autopolyploid sugarcane *Saccharum spontaneum* L. *Accept Nat*
875 *Genet.* 2018;
35
36
37
- 876 48. Paterson AH, Bowers JE, Bruggmann R, Dubchak I, Grimwood J, Gundlach H, et al. The *Sorghum bicolor*
877 genome and the diversification of grasses. *Nature.* 2009;457:551–6.
38
39
40
41
- 878 49. D’Hont A, Ison D, Alix K, Roux C, Glaszmann JC. Determination of basic chromosome numbers in the genus
879 *Saccharum* by physical mapping of ribosomal RNA genes. *Genome.* 1998;41:221–5.
42
43
44
- 880 50. Liu Z, Adams KL. Expression Partitioning between Genes Duplicated by Polyploidy under Abiotic Stress
881 and during Organ Development. *Curr Biol.* 2007;17:1669–74.
45
46
47
48
- 882 51. Ramírez-González RH, Borrill P, Lang D, Harrington SA, Brinton J, Venturini L, et al. The transcriptional
883 landscape of polyploid wheat. *Science.* 2018;361:ear6089.
49
50
51
- 884 52. Zhang Y, Liu Z, Khan AA, Lin Q, Han Y, Mu P, et al. Expression partitioning of homeologs and tandem
885 duplications contribute to salt tolerance in wheat (*Triticum aestivum* L.). *Sci Rep [Internet].* 2016 [cited 2018
886 Aug 16];6. Available from: <http://www.nature.com/articles/srep21476>
52
53
54
55
56
- 887 53. Liu Z, Xin M, Qin J, Peng H, Ni Z, Yao Y, et al. Temporal transcriptome profiling reveals expression
888 partitioning of homeologous genes contributing to heat and drought acclimation in wheat (*Triticum aestivum*
889 L.). *BMC Plant Biol [Internet].* 2015 [cited 2018 Aug 16];15. Available from:
890 <http://www.biomedcentral.com/1471-2229/15/152>
57
58
59
60
61
62
63
64
65

- 891 54. Zou C, Lu C, Shang H, Jing X, Cheng H, Zhang Y, et al. Genome-Wide Analysis of the *Sus* Gene Family in
892 Cotton: Comprehensive Analysis of Cotton *Sus* Genes. *J Integr Plant Biol.* 2013;55:643–53.
- 1
893 55. Rawat R, Schwartz J, Jones MA, Sairanen I, Cheng Y, Andersson CR, et al. REVEILLE1, a Myb-like
894 transcription factor, integrates the circadian clock and auxin pathways. *Proc Natl Acad Sci.* 2009;106:16883–
895 8.
- 6
896 56. Seo PJ, Ryu J, Kang SK, Park C-M. Modulation of sugar metabolism by an INDETERMINATE DOMAIN
897 transcription factor contributes to photoperiodic flowering in Arabidopsis: Sugar and photoperiodic
898 flowering. *Plant J.* 2011;65:418–29.
- 10
899 57. Papini-Terzi FS, Rocha FR, Vêncio RZ, Felix JM, Branco DS, Waclawovsky AJ, et al. Sugarcane genes
900 associated with sucrose content. *BMC Genomics.* 2009;10:120.
- 14
901 58. Kim C, Wang X, Lee T-H, Jakob K, Lee G-J, Paterson AH. Comparative Analysis of Miscanthus and
902 Saccharum Reveals a Shared Whole-Genome Duplication but Different Evolutionary Fates. *Plant Cell.*
903 2014;26:2420–9.
- 18
904 59. Vieira MLC, Almeida CB, Oliveira CA, Tacuatiá LO, Munhoz CF, Cauz-Santos LA, et al. Revisiting Meiosis in
905 Sugarcane: Chromosomal Irregularities and the Prevalence of Bivalent Configurations. *Front Genet* [Internet].
906 2018 [cited 2018 Aug 27];9. Available from:
907 <https://www.frontiersin.org/article/10.3389/fgene.2018.00213/full>
- 24
908 60. Dal-Bianco M, Carneiro MS, Hotta CT, Chapola RG, Hoffmann HP, Garcia AAF, et al. Sugarcane
909 improvement: how far can we go? *Curr Opin Biotechnol.* 2012;23:265–70.
- 28
910 61. Illumina. FastTrack Services Long Reads Pipeline User Guide. 2013.
- 30
911 62. Myers EW. A Whole-Genome Assembly of *Drosophila*. *Science.* 2000;287:2196–204.
- 32
912 63. Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid, finished microbial
913 genome assemblies from long-read SMRT sequencing data. *Nat Methods.* 2013;10:563–9.
- 35
914 64. de Setta N, Monteiro-Vitorello CB, Metcalfe CJ, Cruz GMQ, Del Bem LE, Vicentini R, et al. Building the
915 sugarcane genome for biotechnology and identifying evolutionary trends. *BMC Genomics.* 2014;15:540.
- 39
916 65. Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, et al. Phytosome: a comparative platform
917 for green plant genomics. *Nucleic Acids Res.* 2012;40:D1178–86.
- 42
918 66. Keller O, Kollmar M, Stanke M, Waack S. A novel hybrid gene prediction method employing protein
919 multiple sequence alignments. *Bioinforma Oxf Engl.* 2011;27:757–63.
- 45
920 67. Eddy SR. Accelerated Profile HMM Searches. Pearson WR, editor. *PLoS Comput Biol.* 2011;7:e1002195.
- 48
921 68. Shearman JR, Sonthirod C, Naktang C, Pootakham W, Yoocha T, Sangsrakru D, et al. The two chromosomes
922 of the mitochondrial genome of a sugarcane cultivar: assembly and recombination analysis using long PacBio
923 reads. *Sci Rep* [Internet]. 2016 [cited 2018 Jan 24];6. Available from:
924 <http://www.nature.com/articles/srep31533>
- 53
925 69. Knudsen T, Knudsen B. CLC Genomics Benchwork 6 [Internet]. 2013. Available from:
926 <http://www.clcbio.com>
- 57
927 70. Smit A, Hubley R, Green P. RepeatMasker Open-4.0 [Internet]. Available from:
928 <http://www.repeatmasker.org>
- 60
61
62
63
64
65

- 929 71. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase Update, a database of
930 eukaryotic repetitive elements. *Cytogenet Genome Res.* 2005;110:462–7.
- 1
931 72. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: two open source ab initio eukaryotic
932 gene-finders. *Bioinforma Oxf Engl.* 2004;20:2878–9.
- 4
933 73. Besemer J, Borodovsky M. GeneMark: web software for gene finding in prokaryotes, eukaryotes and
934 viruses. *Nucleic Acids Res.* 2005;33:W451–4.
- 8
935 74. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith Jr RK, Hannick LI, et al. Improving the Arabidopsis
936 genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Res.* 2003;31:5654–66.
- 11
937 75. Slater GSC, Birney E. Automated generation of heuristics for biological sequence comparison. *BMC*
938 *Bioinformatics.* 2005;6:31.
- 13
939 76. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and
940 applications. *BMC Bioinformatics.* 2009;10:421.
- 18
941 77. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic gene structure
942 annotation using EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome Biol.*
943 2008;9:R7.
- 22
944 78. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale protein
945 function classification. *Bioinformatics.* 2014;30:1236–40.
- 26
946 79. Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. Blast2GO: a universal tool for annotation,
947 visualization and analysis in functional genomics research. *Bioinformatics.* 2005;21:3674–6.
- 29
948 80. Han F, Peng Y, Xu L, Xiao P. Identification, characterization, and utilization of single copy genes in 29
949 angiosperm genomes. *BMC Genomics.* 2014;15:504.
- 33
950 81. Kent WJ. BLAT--the BLAST-like alignment tool. *Genome Res.* 2002;12:656–64.
- 35
951 82. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna, Austria; 2014.
952 Available from: <http://www.R-project.org>
- 38
953 83. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids*
954 *Res.* 2004;32:1792–7.
- 42
955 84. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of
956 transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol.* 2013;14:R36.
- 45
957 85. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format
958 and SAMtools. *Bioinforma Oxf Engl.* 2009;25:2078–9.
- 49
959 86. Finn RD, Coghill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, et al. The Pfam protein families database:
960 towards a more sustainable future. *Nucleic Acids Res.* 2016;44:D279–85.
- 52
961 87. Xu Z, Zhang D, Hu J, Zhou X, Ye X, Reichel KL, et al. Comparative genome analysis of lignin biosynthesis
962 gene families across the plant kingdom. *BMC Bioinformatics.* 2009;10 Suppl 1:S3.
- 55
963 88. Röther D, Poppe L, Morlock G, Viergutz S, Rétey J. An active site homology model of phenylalanine
964 ammonia-lyase from *P. crispum*. *Eur J Biochem.* 2002;269:3065–75.

- 965 89. Calabrese JC, Jordan DB, Boodhoo A, Sariaslani S, Vannelli T. Crystal structure of phenylalanine ammonia
966 lyase: Multiple helix dipoles implicated in catalysis. *Biochemistry*. 2004;43:11403–16.
- 1
967 90. Pilbák S, Tomin A, Rétey J, Poppe L. The essential tyrosine-containing loop conformation and the role of
968 the C-terminal multi-helix region in eukaryotic phenylalanine ammonia-lyases. *FEBS J*. 2006;273:1004–19.
- 4
969 91. Thompson JD, Higgins DG, Gibson TJ. CLUSTAL W: improving the sensitivity of progressive multiple
970 sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice.
971 *Nucleic Acids Res*. 1994;22:4673–80.
- 9
972 92. Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis Version 7.0 for Bigger
973 Datasets. *Mol Biol Evol*. 2016;33:1870–4.
- 12
974 93. Zhang Z, Wood WI. A profile hidden Markov model for signal peptides generated by HMMER. *Bioinforma*
975 *Oxf Engl*. 2003;19:307–8.
- 16
976 94. Shahmuradov IA, Umarov RK, Solovyev VV. TSSPlant: a new tool for prediction of plant Pol II promoters.
977 *Nucleic Acids Res*. 2017;gkw1353.
- 19
978 95. Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, et al. MEME Suite: tools for motif discovery
979 and searching. *Nucleic Acids Res*. 2009;37:W202–8.
- 22
980 96. Claeys M, Storms V, Sun H, Michoel T, Marchal K. MotifSuite: workflow for probabilistic motif detection
981 and assessment. *Bioinformatics*. 2012;28:1931–2.
- 26
982 97. Khan A, Fornes O, Stigliani A, Gheorghe M, Castro-Mondragon JA, van der Lee R, et al. JASPAR 2018:
983 update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic*
984 *Acids Res*. 2018;46:D260–6.
- 30
985 98. Russo PST, Ferreira GR, Cardozo LE, Bürger MC, Arias-Carrasco R, Maruyama SR, et al. CEMiTool: a
986 Bioconductor package for performing comprehensive modular co-expression analyses. *BMC Bioinformatics*
987 [Internet]. 2018 [cited 2018 Aug 16];19. Available from:
988 <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-018-2053-1>
- 36
989 99. Alexa A, Rahnenfuhrer J. topGO: Enrichment Analysis for Gene Ontology.
- 39
990 100. Shannon P. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction
991 Networks. *Genome Res*. 2003;13:2498–504.
- 42
992 101. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinforma Oxf*
993 *Engl*. 2010;26:589–95.
- 45
994 102. McCormick RF, Truong SK, Sreedasyam A, Jenkins J, Shu S, Sims D, et al. The Sorghum bicolor reference
995 genome: improved assembly, gene annotations, a transcriptome atlas, and signatures of genome
996 organization. *Plant J Cell Mol Biol*. 2018;93:338–54.
- 50
997 103. Lê S, Josse J, Husson F. FactoMineR : An R Package for Multivariate Analysis. *J Stat Softw* [Internet]. 2008
998 [cited 2017 Nov 30];25. Available from: <http://www.jstatsoft.org/v25/i01/>
- 53
999 104. Josse J, Husson F. missMDA : A Package for Handling Missing Values in Multivariate Data Analysis. *J Stat*
1000 *Softw* [Internet]. 2016 [cited 2017 Nov 30];70. Available from: <http://www.jstatsoft.org/v70/i01/>

1003

1004

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 1 – Genome sequencing: Technology and assembly details and gene prediction features.

	Description	Genomic DNA	BAC clones
Sequencing and assembly data	Sequencing Data	26 Illumina synthetic long-read libraries	Single end Roche 454 of BAC library clones
	Total Sequence	19 Gb	6,6 Gb
	Genome	1.9 x	0.66 x
	Read length Min/Max	1,500 bp / 22,904 bp	8 bp / 2611 bp
	Mean read length	4,930 bp	368.5 bp
	Assembler Software	Celera Assembler (Overlap Graph)	PHRAP/CONSED
	Total reads used in assembly	3,857,849	17,894,306
	Total assembly size	4.26 Gb	49.6 Mb
	Number of unitigs/contigs + singletons	450,609	463
	Contigs Length Min/Max/Mean	1,500 bp / 468,011 bp / 9,452 bp	11,723 bp / 235,533 bp / 107,129 bp
NG50	41,394 bp	109,618 bp	
N50	13,157 bp	N/A	
Gene prediction features	# genes	373,869	3,550
	# transcripts	374,774	-
	# exons	1,035,764	13,132
	Average GC content	43.20%	44.99%
	Average # exons per gene	2.8	3.7
	Average exon size [bp]	291	271.8
	Median exon size [bp]	171	154
	Average intron size [bp]	352.6	539.2
	Median intron size [bp]	132	139
	Average gene size [bp] with UTR	1,437.80	2,429.20
Median gene size [bp] with UTR	806	1,260.50	
Average gene size [bp] without UTR	1,318.80	2,351.30	
Median gene size [bp] without UTR	771	1,199.50	
Average gene density (kb per gene)	11.4	14	

1008

1009 **Figure captions**

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

Fig. 1 – Gene copy number estimation. (A) Distribution of copy counts for putative single copy genes. A total of 1,592 single copy genes from sorghum, rice and *Brachypodium* matched sugarcane predicted genes. More than 99.9% of the aligned single copy genes are present between one and 15 times in the sugarcane gene models. (B) Copy differentiation between sugarcane coding sequences (CDS) and upstream regions, based on pairwise sequence alignment of gene clusters. Genetic dissimilarity increases with increasing distance from the translation start site. (C) Indel length distribution in sugarcane gene copies. Frame preserving indels are more common than frameshifts for this set of genes.

Fig. 2 – Homo (eo)log expression: Frequency of sugarcane genes plotted against the total number of homo(eo)logs per gene and the number of expressed homo(eo)logs per gene. Genes with cDNAs aligned with FPKM > 1 were considered expressed. Plots show sense (A) and antisense (B) transcripts. Reads from Ion PGM Sequencing were used, as strand orientation is maintained [29].

Fig. 3 – Phylogenetic, expression and TFBS categorization of SuSy and PAL gene family. Phylogenetic analysis of (A) sucrose synthase (SuSy) and (B) phenylalanine-ammonia lyase (PAL) genes from SP80-3280, R570, *S. spontaneum*, and sorghum. SuSy sequences from *Saccharum* ssp [35] were also included. Core promoter analysis suggests ScSuSy2 (C) and most ScPAL (D) as TATA-less and TFBS specific for each clade. The three ScPAL genes marked (*) are present in the same contig. Transposable elements (TEs) were identified within 10 kb upstream from the gene. Heatmap analysis of RNA-Seq data [29] shows more pronounced expression in SP80-3280 internodes of ScSuSy1, ScSuSy2, ScSuSy5 and ScPAL from Clade V. RNA-Seq of leaf tissues indicates more pronounced expression of ScPAL from Clades II and III. ScSuSy4 presents high numbers of TFBS and TE and low expression in all samples.

Fig. 4 – SNP variants. Alignment of sugarcane contigs to the genic regions of sorghum chromosomes (chromosome 1 is on top and 10 is at the bottom). X and Y axes indicate physical distance on each chromosome (mega base pairs, Mb) and the number of single nucleotide variants compared to the sorghum reference genome, respectively. Each dot indicates sorghum genes matching two or more sugarcane contigs.

Fig. 5 – Pseudoassembly of contigs. Multiple correspondence analysis (MCA) with hierarchical clustering of the SP80-3280 assembly against the *S. spontaneum* tetraploid AP85-441 homo(eo)log-resolved assembly [14] and the R570 [13] monoplloid genome. A: SP80-3280 contigs best hits against AP85-441 and R579 chromosomes and corresponding size of the preliminary scaffolds; Cluster = hierarchical cluster from the MCA. B and C: Circos plot of the proportion of proteins from SP80-3280 (classified into one of the 6 clusters or as ‘non-clustered’) that align to the AP85-441 and R570 putative chromosomes, respectively.

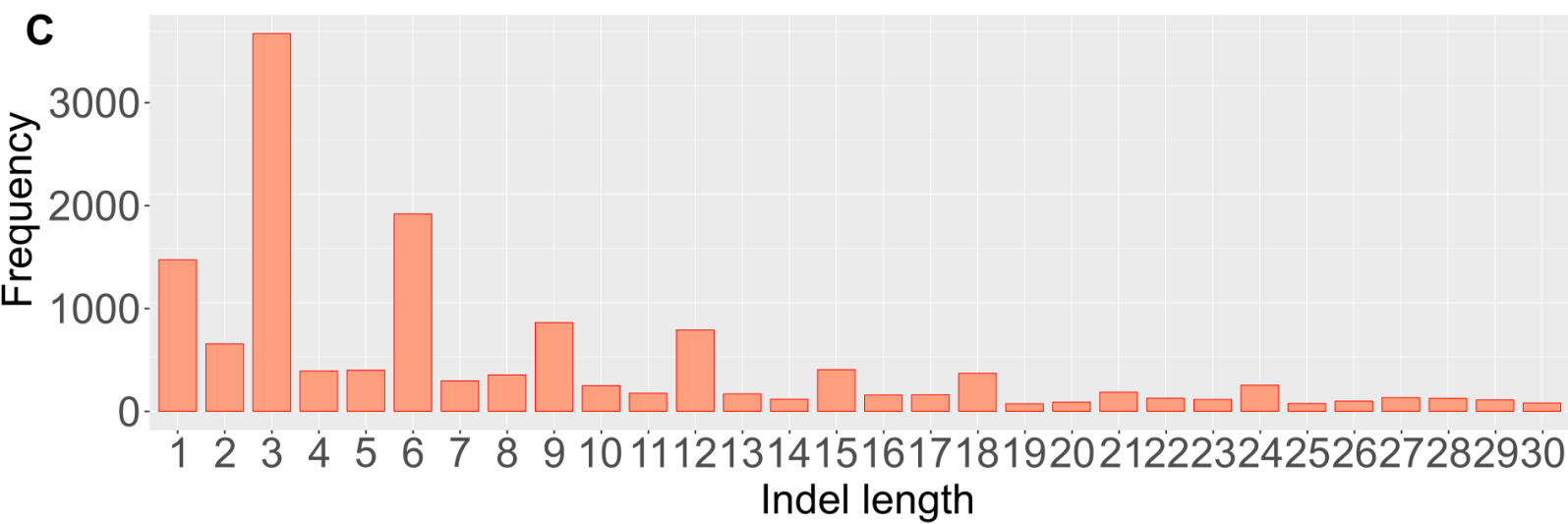
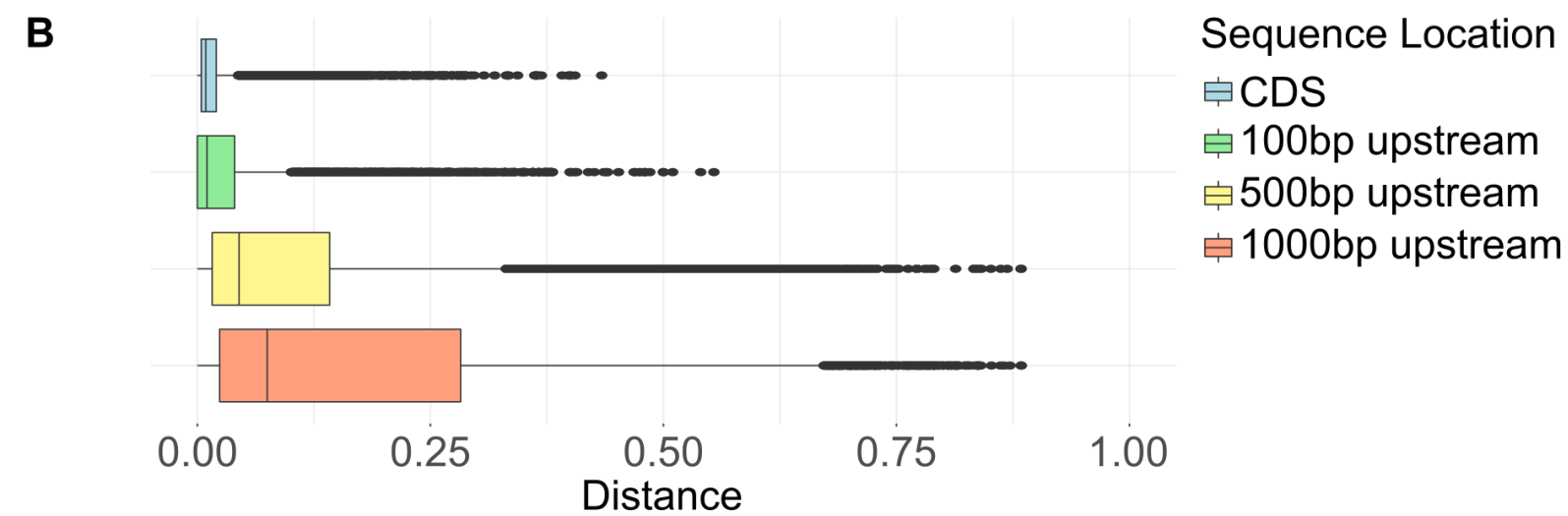
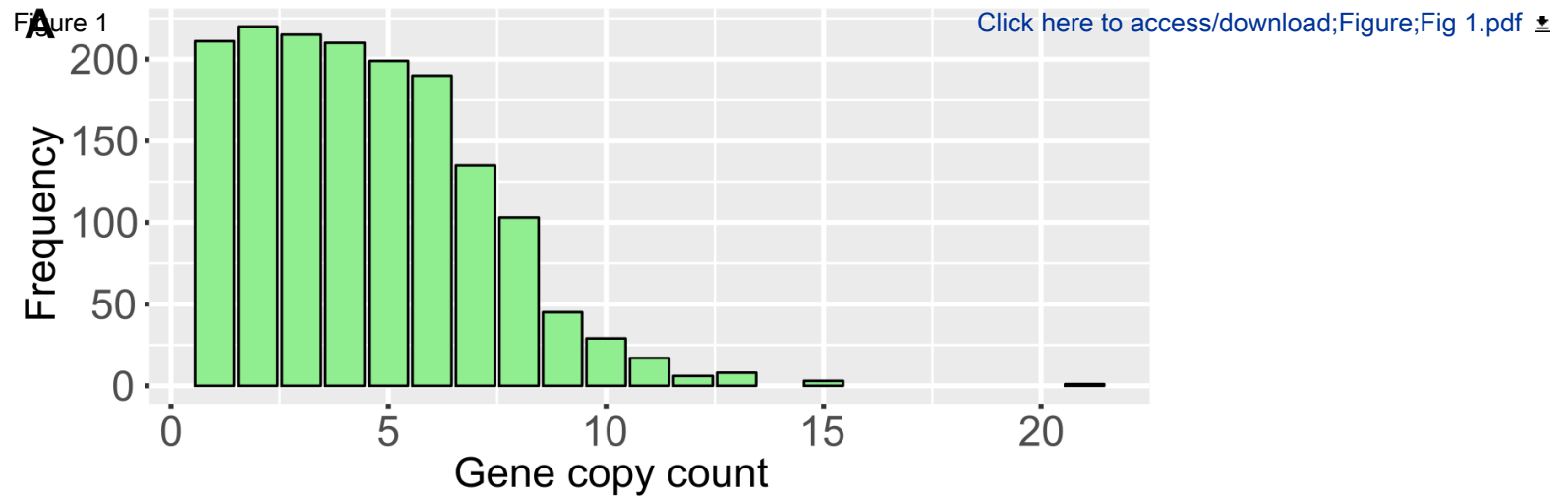
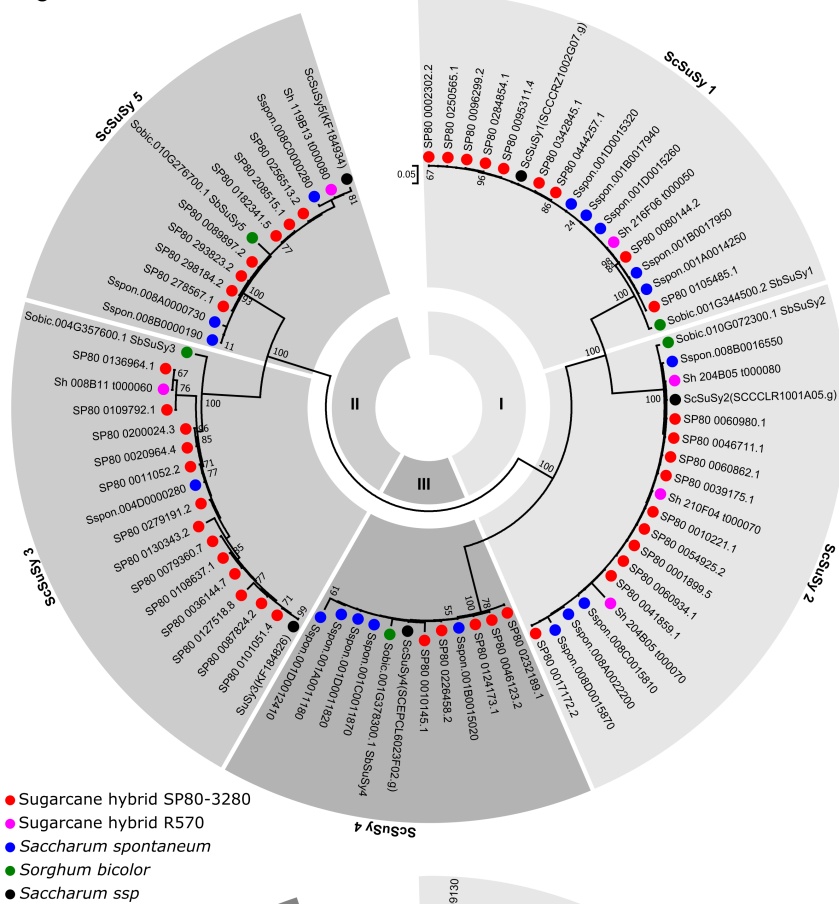
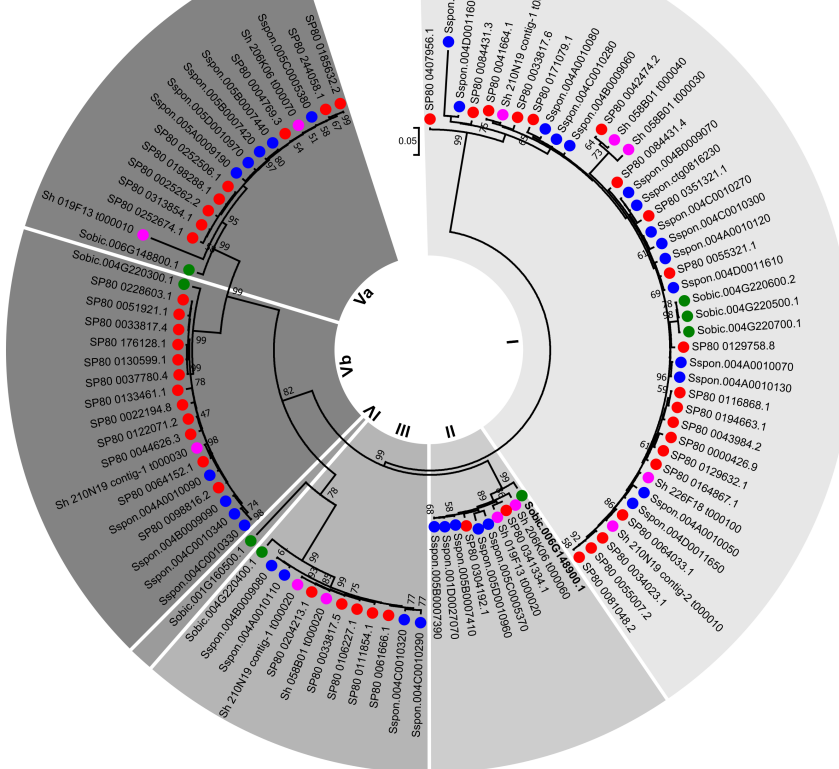


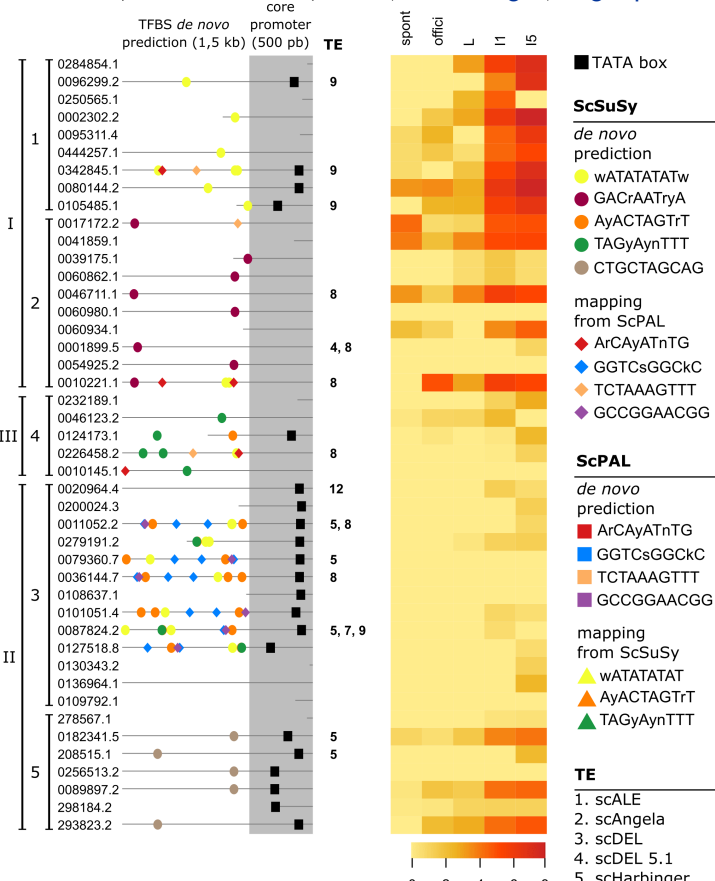
Figure 3



B



C



D

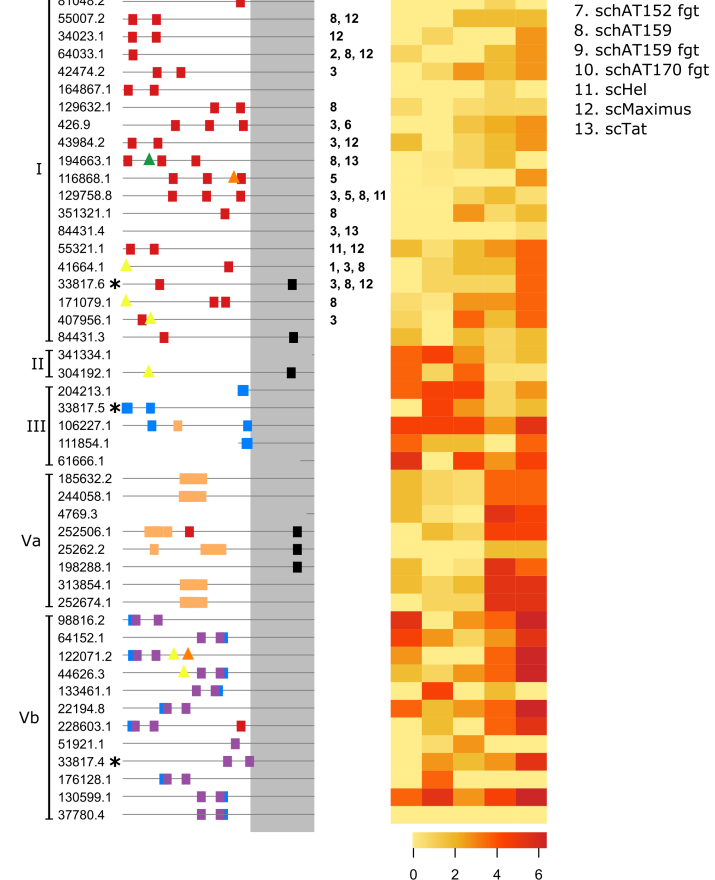
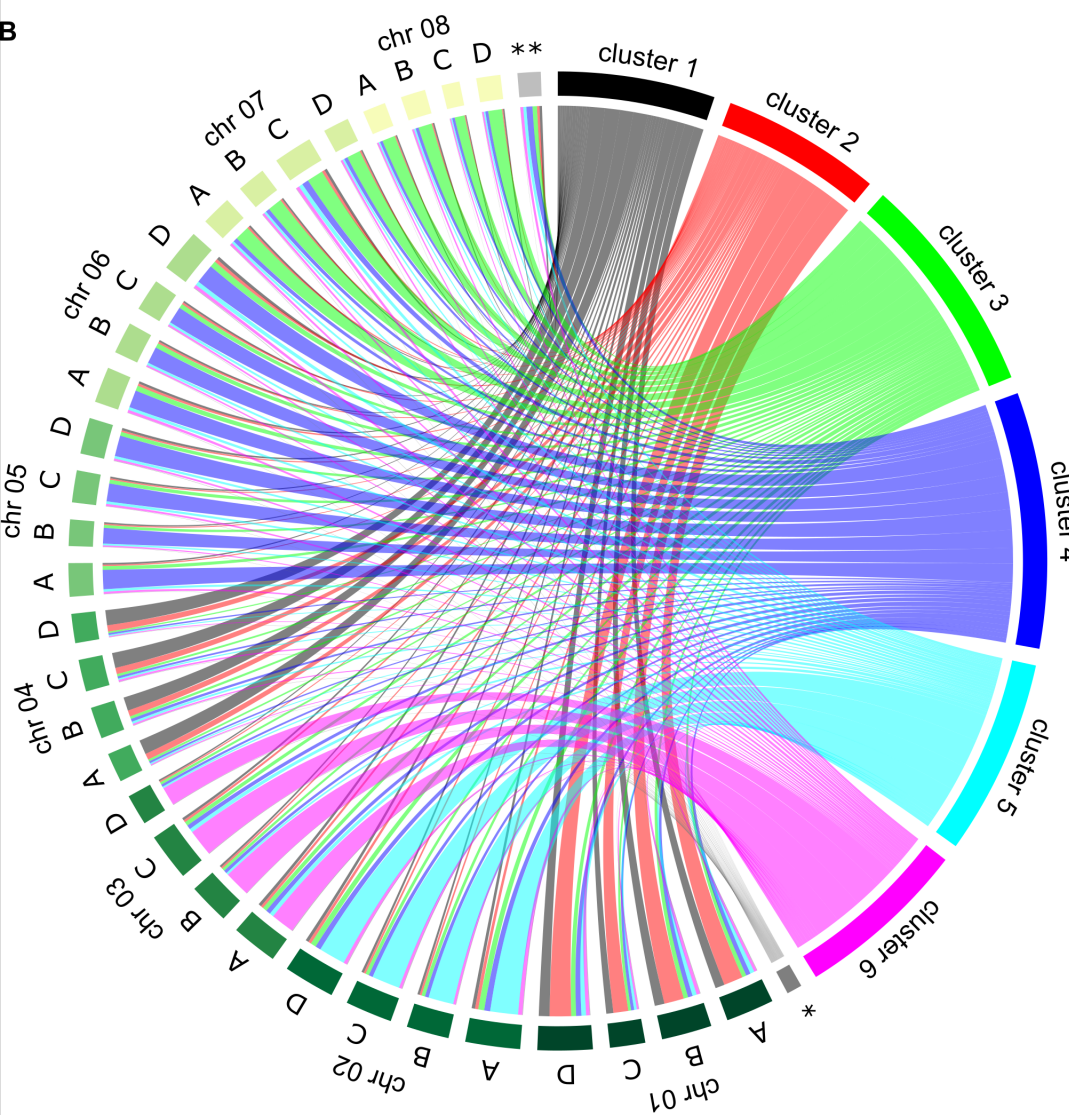


Figure 5

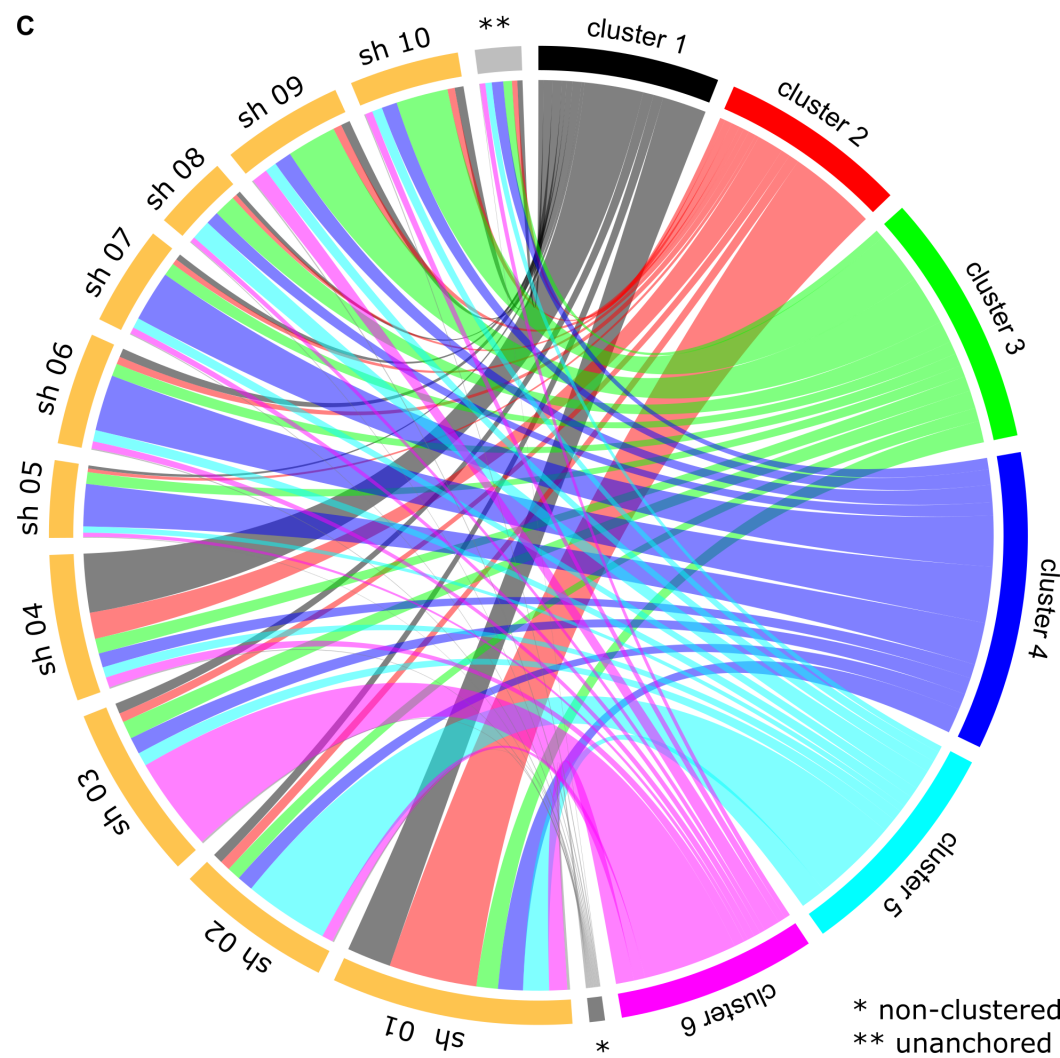
A

Cluster	Number of Contigs	bp	Chromosomal Correspondence	
			<i>S. spontaneum</i>	R570
1	60,150	567,792,642	4	4
2	61,705	574,401,531	1	1
3	87,155	823,254,612	7, 8	8, 9 10
4	90,152	896,362,990	5, 6	5, 6, 7
5	63,996	679,392,733	2	2
6	55,313	565,012,329	3	3
Total	418,471	4,106,216,837	-	-
Original	450,609	4,259,506,050	-	-

B



C





[Click here to access/download](#)

Supplementary Material

[Souza and Van Sluys et al Additional file 1.docx](#)





[Click here to access/download](#)

Supplementary Material

Souza and Van Sluys et al Additional file 2.xls





Universidade de São Paulo
Instituto de Química

Departamento de Bioquímica

Dear Laurie Goodman,
Editor-in-Chief
GigaScience

On behalf of the co-authors, we would like to submit for your consideration the manuscript by Souza, Van Sluys and colleagues entitled “Assembly of the 373K gene space of the polyploid sugarcane genome reveals reservoirs of functional diversity in the world’s leading biomass crop”.

We sequenced the genome of a sugarcane commercial hybrid and were able to produce a high-quality assembly that resolved the homo(eo)logs of the intricate *Saccharum* complex.

Our assembly contains 373,869 genes from a modern commercial cultivar (SP80-3280) used in breeding programs to generate new varieties and that for many years has been the reference of genomic studies. It is important to note that the assembly offers for the first-time access to the gene-space copy-resolved content of interspecific polyploid commercial sugarcane including gene promoters and a gene network analysis. We produced transcriptome data of field grown sugarcane harvested from planting to maturation including four tissues and found regulatory elements of co-expressed genes within a gene network associated with fiber metabolism. They show promoter regions with different transcription factor binding sites associated to differentially expressed genes. In addition, annotation of transposable elements contributed to fine-tuning the observed gene diversity in an otherwise highly syntenic genome with Sorghum.

Our gene space of 373 thousand genes represents a much closer representation of the gene content diversity of *Saccharum* adding valuable new data to previous studies by Gasmour et al., 2018 (doi: 10.1038/s41467-018-05051-5), Zhang et al., 2018 (<https://www.nature.com/articles/s41588-018-0237-2>), Riaño-Pachón and Mattiello, 2017 (doi: <https://doi.org/10.12688/f1000research.11859.2>). In addition, our genome sequence supports its recent allotetraploid nature.



Universidade de São Paulo
Instituto de Química

The data set is a fundamental and large step toward a high-quality chromosome resolved assembly from a current commercial hybrid. We note that as the wheat sequencing genome effort, a multi-initiative is necessary to approach complex genomes, but different from wheat, sugarcane genetics is hindered by it being mainly propagated through cuttings and crosses being performed only for breeding purposes. The availability of a polyploid gene-space will be a valuable step towards breeding in a plant extraordinarily difficult to yield genetic maps.

This Whole Genome Shotgun project is publicly available at NCBI under GenBank Bioproject PRJNA431722. Finally, we assure that all gene/protein names and symbols used in this manuscript are in accordance to approved nomenclature guidelines for Saccharum species.

We hope you find this achievement to be of interest to GigaScience and look forward to hearing from you.

Sincerely,

Glaucia Mendes Souza

Full Professor

Institute of Chemistry

University of São Paulo

Marie-Anne Van Sluys

Full Professor

Biosciences Institute

University of São Paulo