

1 **Genome and transcriptome architecture of allopolyploid**
2 **okra (*Abelmoschus esculentus*)**

3 **Ronald Nieuwenhuis¹, Thamara Hesselink¹, Hetty C. van den Broeck¹, Jan**
4 **Cordewener¹, Elio Schijlen¹, Linda Bakker¹, Sara Diaz Trivino¹, Darush Struss²,**
5 **Simon-Jan de Hoop², Hans de Jong³ and Sander A. Peters*¹.**

6 ¹Business Unit of Bioscience, cluster Applied Bioinformatics, Wageningen University and
7 Research, Droevendaalsesteeg 1, 6708 PB Wageningen, The Netherlands.

8 ²East-West International B.V., Heiligeweg 18, 1601 PN Enkhuizen, the Netherlands.

9 ³Laboratory of Genetics, Wageningen University, Droevendaalsesteeg 1, 6708 PB,
10 Wageningen, The Netherlands.

11 *corresponding author; Sander A. Peters; tel: +31-317-481123; e-mail:
12 sander.peters@wur.nl

13 **Running title:** The okra genome and transcriptome profile

14 **Key words:** Okra, allopolyploid, *Malvaceae*, chromosomes, genome, transcriptome,
15 annotation, BUSCO, telomere, rRNA genes, polyphenols, flavonoid biosynthesis

16

17

18

19

20

21

22

23

24

1 **Abstract**

2 We present the first annotated genome assembly of the allopolyploid okra (*Abelmoschus*
3 *esculentus*). Analysis of telomeric repeats and gene rich regions suggested we obtained whole
4 chromosome and chromosomal arm scaffolds. Besides long distal blocks we also detected short
5 interstitial TTTAGGG telomeric repeats, possibly representing hallmarks of chromosomal speciation upon
6 polyploidization of okra. Ribosomal RNA genes are organized in 5S clusters separated from the 18S-5.8S-
7 28S units, clearly indicating an S-type rRNA gene arrangement. The assembly is consistent with
8 cytogenetic and cytometry observations, identifying 65 chromosomes and 1.45Gb of expected genome
9 size in a haploid sibling. Approximately 57% of the genome consists of repetitive sequence. BUSCO
10 scores and A50 plot statistics indicated a nearly complete genome. Kmer distribution analysis suggests
11 that approximately 75% has a diploid nature, and at least 15% of the genome is heterozygous. We did
12 not observe aberrant meiotic configurations, suggesting there is no recombination among the sub-
13 genomes. BUSCO configurations as well as k-mer clustering analysis pointed to the presence of at least
14 2 sub-genomes. These observations are indicative for an allopolyploid nature of the okra genome.
15 Structural annotation, using gene models derived from mapped IsoSeq transcriptome data, generated
16 over 130,000 putative genes. Mapped transcriptome data from public okra accessions of Asian origin
17 confirmed the predicted genes, showing limited genetic diversity of 1SNP/2.1kb. The discovered genes
18 appeared to be located predominantly at the distal ends of scaffolds, gradually decreasing in abundance
19 toward more centrally positioned scaffold domains. In contrast, LTR retrotransposons were more
20 abundant in centrally located scaffold domains, while less frequently represented in the distal ends. This
21 gene and LTR-retrotransposon distribution is consistent with the observed heterochromatin organization
22 of pericentromeric heterochromatin and distal euchromatin. The derived amino acid queries of putative
23 genes were subsequently used for phenol biosynthesis pathway annotation in okra. Comparison against
24 manually curated reference KEGG pathways from related *Malvaceae* species revealed the genetic basis
25 for putative enzyme coding genes that likely enable metabolic reactions involved in the biosynthesis of
26 dietary and therapeutic compounds in okra.

27

28

29

30

31

32

1 Introduction

2 The well-known okra (*Abelmoschus esculentus*) vegetable belongs to the family *Malvaceae*,
3 comprising more than 244 genera and over 4,200 species. The *Malvaceae* are divided into 9 subfamilies
4 of which okra belongs to the subfamily *Malvoideae*. *Abelmoschus* is closely related to *Hibiscus* species
5 like *Hibiscus rosa-chinensis* or Chinese rose and *Hibiscus cannabinus* or Kenaf, which is exemplified by
6 the beautiful characteristic Hibiscus-like flowers that both genera display. Based on genetic differences,
7 *Abelmoschus* has now been placed in a separate genus from *Hibiscus* though. Okra is flowering
8 continuously and is self-compatible, however cross-pollination up to 20% has been reported. Its
9 characteristic hermaphroditic flowers usually have white or yellow perianths, consisting of five petals and
10 five sepals, whereas calyx, corolla and stamens are fused at the base. Other well-known species in the
11 *Malvaceae* are cocoa (*Theobroma cacao*), cotton (*Gossypium hirsutum*), and *Tilia* species like lime tree,
12 the mangrove *Heritiera* species and durian (*Durio zibethinus*). The genus *Abelmoschus* contains 11
13 species, four subspecies and five varieties (Li *et al.*, 2020) of which most members have economic value.
14 Okra or 'lady's finger' is a low-calorie vegetable, mainly cultivated for its fruits that are harvested while
15 still unripe, containing a large variety of nutrients and elements essential for daily human consumption,
16 such as vitamins, flavonoids, minerals, and other health components such as folate and fibers (Muimba-
17 Kankolonga, 2018; Wu *et al.*, 2020). For example, total polyphenol extracts from okra fruits, containing
18 flavonoids such as myricitin and quercitin, have been demonstrated for their antidiabetic activity in obese
19 rats suffering from type 2 diabetes mellitus (Peter *et al.*, 2021). These health compounds and additional
20 nutritional qualities make okra an appreciated vegetable in many parts of the tropics and subtropics of
21 Asia, Africa and America, gaining rapidly in popularity. Global production has increased yearly since
22 1994, reaching 10M tonnes in 2019 and covering some 2.5M ha (<http://faostat.fao.org>), with Asia having
23 the largest production share of almost 70%, of which India alone is currently annually producing more
24 than 4M tonnes. However, its production is challenged by a range of pathogens and insect pests, such as
25 powdery mildew and blackmold (*Cerospora abelmoschii*), bacterial blight disease (*Xanthomonas*
26 *campestris* p.v. *malvacearum*), mycoplasmas, nematodes, worms and insects such as whitefly (*Bemisia*
27 *tabaci*), thrips (*Thrips palmi*), cotton leafhopper (*Amarasca biguttula*) and aphids (*Aphis gossypii*).
28 Besides feeding damage, these vectors can transmit viruses such as Yellow Vein Mosaic Virus (YVMV), a
29 geminivirus, causing crop losses of up to 80-90% without pest control (Benchasri, 2012; Muimba-
30 Kankolonga, 2018; Dhankhar and Koundinya, 2020; Lata *et al.*, 2021). Typical symptoms of YVMV
31 infected okra plants are a stunted growth, with vein and veinlets turning yellow in colour, producing seed
32 pods that are small, distorted, and chlorotic. Crop loss may be reduced to 20-30%, by controlling insect
33 pests with rather harmful and toxic chemicals and insecticides (Ali *et al.*, 2005), causing considerable

1 collateral damage to the ecosystem. Moreover, increased insect tolerance to pesticides has led to over-
2 use and mis-use of chemicals, leaving unhealthy high levels of pesticide residues (Benchasri, 2012).
3 Although there are some YVMV tolerant okra genotypes, such as Nun1144 and Nun1145
4 (Venkataravanappa *et al.*, 2013), the genetic basis for this tolerance has not been identified. Besides a
5 need for disease resistance, other breeding challenges and demands include maximizing production,
6 unravelling the genetic basis for abiotic stress tolerance, and the need to develop double haploid lines
7 enabling the study of recessive gene traits (Dhankhar and Koundinya, 2020).

8 To meet current demands and challenges, accelerated breeding is urgently needed. Presently,
9 several methods of breeding for improvement in okra are being used, such as pure line selection,
10 pedigree breeding, as well as mutation and heterosis breeding (Dhankhar and Koundinya, 2020). These
11 methods are very time consuming though, and often involve laborious analyses over multiple
12 generations. Despite wide genetic variation available among wild relatives of okra, significant crop
13 improvement by introgression breeding, has not been achieved due to hybridization barriers. Advanced
14 breeding is further hampered due to the lack of sufficient molecular markers, linkage maps and reference
15 genome, and this in turn has impeded genome and transcriptome studies. Molecular studies have further
16 been complicated due to the presence of large amounts of mucilaginous and polyphenolic compounds in
17 different tissues, interfering with the preparation of genetic materials (Takakura and Nishio, 2012; Lata
18 *et al.*, 2021). Furthermore, correct *de novo* assembly is presumed to be complex because of the
19 expected large genome and transcriptome size and the highly polyploid nature of the genome. Salameh
20 (2014) reported flow-cytometric estimates of nuclear DNA size estimations with 2C values ranging from
21 3.98 to 17.67 pg, equaling to genome sizes between 3.8 to 17.3 Gbp. In addition, chromosome counts
22 demonstrated a huge variation, ranging from 2n=62 to 2n=144, with 2n=130 as the most frequently
23 observed chromosome number (Benchasri, 2012; Merita *et al.*, 2012). These findings have led to further
24 assumptions on the geographical origin of cultivated *A. esculentus*, speculating that a 2n=58 specie *A.*
25 *tuberculatus* native from Northern India and a 2n=72 specie *A. ficulneus* from East Africa might have
26 hybridized followed by a chromosome doubling, giving rise to an allopolyploid *Abelmoschus* hybrid with
27 2n=130 (Joshi and Hardas, 1956; Siemonsma, 1982; Benchasri, 2012; Merita *et al.*, 2012). However,
28 genomic, genetic and cytological information is scanty, limiting the possibilities to further understand the
29 hereditary constituent of the crop. In this study we benefitted from naturally occurring okra haploids,
30 circumventing heterozygosity in the reconstruction of composite genome sequences, supporting faithful
31 genome reconstruction (Langley *et al.*, 2011). Here we present a detailed insight into the complex
32 genome and transcriptome architecture of an okra cultivar and its haploid descendent, using cytogenetic
33 characterization of its mitotic cell complements and meiosis, and advanced sequencing and assembly

1 technologies of the haploid genome, providing basic scientific knowledge for further evolutionary studies
2 and representing a necessary resource for future molecular based okra breeding. Furthermore, we
3 provide a structural and functional genome annotation that is of paramount importance to understand
4 plant metabolism (Weissenborn *et al.*, 2017) and the genetic basis for the enzyme coding genes that
5 enable metabolic reactions involved in the biosynthesis of dietary and therapeutic compounds in okra.

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

1 **Results and discussion**

2 *Cytogenetic characterization of the okra crop*

3 As okra is known to contain large numbers of chromosomes that differ between cultivars, we
4 first established chromosome counts and morphology in the cultivar used in this study. Actively growing
5 root tips were fixed and prepared for cell spread preparations following a standard pectolytic enzyme
6 digestion and air-drying protocol, and DAPI fluorescence microscopy (Kantama *et al.*, 2017). In this
7 diploid red petiole phenotype plants, we counted 130 chromosomes in late prophase and metaphase cell
8 complements (Figures 1a). Chromosomes measured 1-2 μm , often show telomere to telomere
9 interconnections (Figure 1b) and were clearly monocentric (Figure 1a, arrows). In addition, a few
10 chromosomes displayed a less condensed distal region at one of the chromosome arms and satellites
11 (Figure 1c,d), which we interpreted as the nucleolar organiser region (NOR) of the satellite chromosome.
12 Interphase nuclei showed well differentiated heterochromatic domains or chromocenters, most of them
13 with more than 130 spots, although a small number of nuclei decondensed most of its heterochromatin,
14 leaving only a striking pattern of about 10 chromocenters (Figure 1e). We next applied flow cytometry on
15 DAPI stained nuclei, using young leaf material from five normally growing red petiole phenotype
16 plantlets. Since we expected a considerable DNA content for okra nuclei, we decided to use a reference
17 sample from Agave (*Agave americana*), which has a known DNA content of 15.90 picogram. Surprisingly,
18 in comparison to the Agave reference flow cytometric profile, the 2C DNA amount for the normal okra
19 plant was estimated at 2.99 pg \pm 0.01 (Table S1). This amount is equivalent to a genome size of
20 approximately 2.92 Mbp. In contrast to 130 chromosomes for the diploid okra, we observed a
21 chromosome number of 65, which was considered a haploid (Figure 1f). The genome size for this haploid
22 okra was estimated at 1.45Mbp. This plant was feeble, lagging in growth and unfortunately died
23 precociously, but encouraged us to seek for haploid offspring in later samples of reared young okra
24 plants. Such natural haploids, of which a dwarf form of cotton (*Gossypium*) was discovered in 1920 as
25 the first haploid angiosperm with half the normal chromosome complement (Dunwell, 2010), are
26 assumed to result from asexual egg cell (gynogenetic) reproduction (Noumova, 2008). We took
27 advantage of the fact that the diploid hybrid cultivar has a green recessive petiole female parent and a
28 red dominant petiole male (Figure S1) (Portemer *et al.*, 2015). Offspring with the green petiole trait
29 lacks the dominant paternal allele and hence can be used as a diagnostic marker for identifying haploid
30 offspring. Accordingly, we selected additional haploid offspring of which one plant was used for
31 sequencing.

1 For the analysis of homologous pairing, chiasma formation and chromosome segregation in
2 diploid okra plants we studied male meiosis in pollen mother cells from young anthers (Kantama *et al.*,
3 2017). Pollen mother cells at meiotic stages are filled with fluorescing granular particles in the
4 cytoplasm, which makes it notoriously difficult to see fine details in chromosome morphology. By long
5 enzymatic digestion and acetic acid maturation we still could make the following details visible:
6 pachytene was strikingly diploid-like with clear bivalents showing denser pericentromere regions and
7 weaker fluorescing euchromatin distal parts (Figure 1f). We did not observe clear inversion loops
8 indicative for inversion heterozygosity or pairing partner switches that demonstrate homoeologous
9 multivalents or heterozygous translocation complexes, however the occurrence for such chromosome
10 structure variants could not be excluded. Cell complements at diakinesis displayed that most (if not all)
11 chromosome configurations were bivalents, supporting a diploid like meiosis (Figure 1g). We did not see
12 univalents or laggards at later stages, and pollen were strikingly uniform and well stained (data not
13 shown).

14 *Okra haploid genome reconstruction*

15 Based on public reports (Benchasri, 2012; Salameh, 2014) and cytological analysis presented
16 above, we applied several technologies for genome reconstruction of the okra haploid individual. 10X
17 Genomics linked read information was used to obtain sequence information in the 100kb to 150kb range.
18 This microfluidics-based technology combines barcoded short-read Illumina sequencing, allowing a set of
19 150 bp paired-end reads to be assigned to large insert molecules. We produced 800 Gbp of linked read
20 sequencing data from three libraries with an average GC-content of 34% (Table S2). Furthermore, we
21 applied PacBio Circular Consensus Sequencing (CCS), generating 1,400 Gbp of polymerase reads of up to
22 150kb from circularized insert molecules from three sequence libraries with fragment insert sizes of 10,
23 14 and 18 Kbp respectively. Polymerase reads were subsequently processed into consensus or so-called
24 HiFi reads with an average sub-read length of approximately 13.6 kb and a sequence error rate less than
25 1% (Table S2). Over 93% of 1,000 randomly sampled CCS reads had a best BlastN hit to species from
26 the *Malvaceae* family with *Gossypium* ranking first in number of hits, indicating the consistent taxonomic
27 origin, in contrast to the *Abelmoschus* species that are apparently less represented in the NCBI sequence
28 database (Table S3). Furthermore, the organellar DNA content was sufficiently low (Table S4),
29 illustrating the efficiency of our nuclear DNA sample preparations. Upon assembling the HiFi reads with
30 the Hifiasm assembler (Cheng *et al.*, 2021), we obtained 3,051 high quality primary contigs with an N50
31 contig length of 18.9 Mb (Table S5). The incremental sequence assembly displayed in the A50 plot
32 (Figure S2) shows a plateau genome size of approximately 1.35 Gbp, which agrees with the nuclear 2C
33 DNA content. The assembly also resulted in 972 alternative contigs, although their total size of 31 Mbp

1 was small, indicating a highly consistent primary assembly. We nevertheless assessed the origin of
2 alternative contigs, using a taxon annotated GC (TAGC) screen (Kumar *et al.*, 2013), providing a means
3 to discriminate between on-target and off-target genomic sequence based on the combined GC content
4 and read coverage and corresponding best matching sequence in annotated databases. The distribution
5 and specific classes of Blast hits indicated that approximately 30% of the alternative contigs could be
6 mapped against annotated sequences (Figure S3), while two thirds were of unknown origin. Alternative
7 contigs had a GC content of 47.6%, which was proportionally higher compared to primary contigs.
8 Furthermore, BlastN hits pointed to a fungal and, to a lesser extent, a bacterial origin. Thus, the smaller
9 sized alternative contigs represented yet a minor contamination in the gDNA sample.

10 We next physically mapped the genome with BioNano Genomics technology to determine the
11 genome structural organization. We produced 4.88 Tb of unfiltered genome map data with an N50
12 molecule size of 90.18 kb (Table S2) and a label density of 15.9 per 100kb (Table S6) from nuclei
13 preparations of leaf samples. Size filtering for molecules larger than 100kb left approximately 1.2 Tb of
14 genome mapping data with an N50 molecule size of 206.8 kb (Table S6). Next, molecules, having
15 matching label position and distance, were *de novo* assembled into 216 genome maps with an N50
16 length of 12.98 Mbp and a total length of 1248.8 Mbp, representing an effective coverage of 375X (Table
17 S6). The genome map size was consistent with the genome sequence assembly size of 1.2 Gbp and thus
18 provided high quality ultra-long-range information for further scaffolding. For that, genome maps were
19 aligned with the *in silico* DLE restriction maps from primary sequence contigs and assembled into higher
20 order scaffolds. The alignment required the cutting of 1 optical map and 4 sequence assembly contigs to
21 resolve conflicts between Bionano maps and sequence contigs respectively, indicating a consistent
22 orientation and order between both. The resulting hybrid assembly was substantially less fragmented,
23 yielding 80 scaffolds with an N50 scaffold size of 18.93 Mb and a total length of 1.19 Gbp, of which the
24 largest scaffold sized more than 29 Mbp (Table 1). Additional scaffolding with 10X Genomics linked reads
25 finally yielded 78 scaffolds (Table 1). Approximately 57% of the individual molecules could be mapped
26 back to the hybrid assembly, suggesting a high confidence genome scaffold.

27 *BUSCO analysis and topology of orthologs*

28 To assess the completeness of the genome assembly we screened for BUSCO gene
29 presence/absence (Simaõ, 2015) with 2,326 reference orthologs from the eudicots_odb10 dataset.
30 Based on a best tBlastN hit, 2270 (98%) core genes in 78 scaffolds were detected as 'Complete'
31 orthologs (Table 2). Of these, 284 (12.2%) genes were detected as a single copy ortholog. A very small
32 amount (0.3%) was classified as 'Fragmented', whereas 32 core genes (1.3%) could not be detected,

1 classifying them as 'Missing'. These missing BUSCO genes were confirmed to be missing in the
2 alternative contigs as well. We further grouped 2004 (86.2%) multiplied ortholog genes according to
3 their copy number. A majority of 1150 (49.4%) and 843 (36.2%) orthologs were detected as duplicated
4 and triplicated genes respectively. Interestingly, we found seven and three core genes that were
5 quadruplicated (0.7%) and quintuplicated (0.1%) respectively, and detected one septuplicated core
6 gene, pointing to a complex polyploid nature of the okra genome. To get more insight into the sub-
7 genome organization, the genomic position and topology of ortholog gene copies was assessed. This
8 revealed duplicated BUSCO genes predominantly occurring on two contigs, whereas only a single
9 duplicated ortholog was detected on one contig (Table S7). Both tandem copies were spaced within 1 kb,
10 thus likely representing paralogous genes. Out of 800 triplicate BUSCO's, 794 (99%) occurred on three
11 contigs, representing three alleles of the same core gene, whereas only six sets (1%) of triplicate core
12 genes were positioned on two contigs. Also, quadruplicate, quintuplicate and septuplicate BUSCO's
13 mainly occurred on three contigs. The ortholog copies of these groups manifested in a tandem
14 configuration, probably also representing paralogs. Tandemly arranged ortholog copies on the same
15 contigs always showed less sequence distance than between ortholog copies on different contigs.
16 Moreover, the ortholog copies of the septuplicate core gene were dispersed over three contigs. Of these,
17 one contig displayed a triplet, whereas the two other contigs each contained gene copies in a doublet
18 configuration. The triplet consisted of two closely related and one more distantly related paralog. The
19 observed distribution of BUSCO orthologs thus pointed to at least two sub-genomes. However, at this
20 point we could not rule out a higher number of sub-genomes, which might not be discriminated because
21 of a low allelic diversity. Considering a sub-genomic organization for okra, we presume that 284 'single
22 copy' BUSCO genes are either truly unique, or they are maintained as gene copies with indistinguishable
23 alleles.

24 Several examples of BUSCO duplication levels in homozygous and heterozygous diploids as well
25 as in auto and allopolyploids have previously been presented for other species. For example in allotetraploid
26 ($2n=4x=38$) *Brassica napus* 90% of BUSCO's are duplicated, whereas only 14.7% in its diploid relative
27 *Brassica campestris* ($2n=2x=18$) or Chinese cabbage is duplicated (Table S8). The allotetraploid white
28 clover ($2n=4x=32$) (*Trifolium repens*), that has suggested to be evolved from 2 related diploid species *T.*
29 *occidentale* ($2n=2x=16$) and *T. pallescens* ($2n=2x=16$), has 57% of duplicated BUSCO's, compared to
30 10% and 11% of BUSCO duplicates in its diploid ancestral relatives respectively (Griffiths *et al.*, 2019).
31 Duplicated BUSCO's in hexaploid bamboo *B. amplixicaulis* ($2n=6x=72$) has increased to 57% compared
32 to 35% in its diploid bamboo relative *O. latifolia* ($2n=2x=22$) (Guo *et al.*, 2019). Significant differences
33 were also observed in BUSCO scores between heterozygous and homozygous diploid *Solanaceae*. For

1 example the heterozygous *S. tuberosum* RH potato ($2n=2x=24$) showed 74.1% of its BUSCO's
2 duplicated, which was significantly more than 4.3% of duplicated BUSCO's detected in diploid inbred *S.*
3 *chacoense* M6 potato ($2n=2x=24$), and 9.5% detected in autotetraploid inbred *S. tuberosum* potato
4 (Kyriakidou *et al.*, 2020). Considering these trends, the BUSCO copy numbers in the okra genome likely
5 point to an allopolyploid nature. Our results confirm the allopolyploid nature of okra as reported by Joshi
6 and Hardas (1956).

7 *K-mer counts and smudgeplot analysis*

8 To estimate heterozygosity level, repetitiveness, genome size and ploidy levels, we determined
9 kmer counts from raw Illumina and HiFi reads. We compared the 21-mers counts for okra to two related
10 allotetraploid cotton species (*G. barbadense* and *G. hirsutum*), each having a genome of approximately
11 2.3Gb, and subsequently visualized the readout with SMUDGEPLOT (Ranallo-Benavidez *et al.*, 2020)
12 (Figure S4). The k-mer based genome size estimation for the haploid okra amounted to 1.2 Gbp,
13 approximating the NGS assembly size. Approximately 75% of the okra kmers was assigned to an 'AB'
14 type (Figure S4). Thus, a major part of the okra genome apparently behaved as a diploid, which is
15 consistent with our cytological observations of a diploid like meiosis, and also coincides with the high
16 number of duplicated BUSCO scores. Approximately 15% of all kmer pairs showed a triploid behaviour
17 ('AAB-type'). Furthermore, the 'AAAB'-kmer type seemed more prominent than the 'AABB'-kmer type.
18 Previously, published kmer readouts for *G. barbadense* and *G. hirsutum* showed that at least 50% of the
19 cotton genomes behaved like a diploid, almost a quarter displayed a triploid behaviour and 14% of the
20 kmers showed tetraploid characteristics. Furthermore, cotton kmer distributions showed the 'AAAB' type
21 more frequently occurring than the 'AABB' kmer type, which was suggested to be a characteristic for
22 allopolyploids (Ranallo-Benavidez *et al.*, 2020). Thus GENOMESCOPE and SMUDGEPLOT readouts for okra point to
23 an allopolyploid nature of the genome, though less complex and smaller sized than anticipated.

24 *Transcriptome profiling and structural annotation*

25 In addition to sequencing the nuclear genome, we generated approximately 1.2 Tbp of IsoSeq
26 data from multiple tissues including leaf, flower buds and immature fruits to profile the okra
27 transcriptome. The polymerase mean read lengths of up to 86kb benefitted the processing into high
28 quality CCS reads with a mean length of 4.7 kb ($\sigma=378$ bp), indicating the efficient full length transcript
29 sequencing (Table S8). The CCS reads were used as transcript evidence for okra gene modelling with the
30 AUGUSTUS and GENEMARK algorithms from BRAKER2. We subsequently annotated the 78 largest okra
31 scaffolds with 130,324 genes. Predicted genes had an average length of 2537 nts, whereas the average
32 per gene intron and coding sequence lengths amounted to 307 nts and 2497 nts respectively (Table 3).

1 Coding regions showed low sequence diversity, as only 1109 and 8127 SNPs could be called from full-
2 length transcripts of the okra haploid and an unrelated diploid individual respectively (Table S9).
3 Strikingly, the discovered genes appeared to be predominantly located at the distal ends of scaffolds,
4 gradually decreasing in abundance toward more centrally positioned scaffold domains. In contrast, LTR
5 retrotransposons were more abundant in centrally located scaffold domains, while less frequently
6 represented in the distal ends (Figure 2). A comparable distribution of gene and LTR-retrotransposon
7 regions has been observed for other species such as tomato. The gene and LTR-retrotransposon predicts
8 a heterochromatin organization of pericentromere heterochromatin and distal euchromatin as shown in
9 Figure 1g and inset. This pattern is common in species with small or moderate chromosome size like
10 Arabidopsis, rice and tomato. Apparently, okra also has relatively small sized chromosomes as is
11 substantiated by our cytological observations. The gene-rich regions predominantly occur in euchromatin
12 rich distal chromosome ends and gradually decrease towards the repeat rich more condensed
13 pericentromeric heterochromatin, whereas LTR-retrotransposons were more frequently distributed in
14 pericentromeric heterochromatin (Peters *et al.*, 2009; Tomato Sequencing Consortium, 2012; Aflitos *et*
15 *al.*, 2014). Our observations thus suggest a similar chromatin architecture for okra chromosomes.
16 Approximately 51% of the assembled genome was found in the repetitive fraction with 20.26% of
17 repeats unclassified (Table 3). A substantial part (28.69%) consisted of retroelements, of which 24.8%
18 and 1.17% was identified as retrotransposon and DNA transposon respectively. *Gypsy* and *Ty1/Copia*
19 retroelements, spanning 15.26% and 9.02% of the assembled genome respectively, appeared to be
20 most abundant (Table 3).

21 The repeat screening also revealed stretches of the Arabidopsis telomere TTTAGGG motif, flanking gene
22 rich regions at distal scaffold ends (Figure 3). In plants such repeats usually occur in high copy numbers
23 at the distal ends of chromosomes, constituting telomeres that protect the terminal chromosomal DNA
24 regions from progressive degradation and preventing the cellular DNA repair mechanism from mistaking
25 the ends of chromosomes for a double stranded break. Indeed, we found blocks of TTTAGGG units in
26 high copy numbers positioned at both ends for 49 scaffolds, whereas 25 scaffolds had a telomere repeat
27 block at one end, in total 123 telomeres at the end of 130 chromosome arms. This repeat distribution
28 suggested full length chromosome scaffolds and capturing the majority of 65 chromosome ends of the
29 haploid okra genome and again confirms the relatively small sized okra chromosomes. Besides long distal
30 blocks we also detected short interstitial TTTAGGG blocks. These interspersed non-telomeric short
31 TTTAGGG repeats possibly reflect footprints of internalized telomeres that may have arisen from end-to-
32 end fusion of chromosomes (Baird, 2017), possibly representing hallmarks of chromosomal speciation
33 upon allopolyploidization of okra. Another substantial fraction of repeats originated from ribosomal

1 genes. BlastN analysis, using *Gossypium hirsutum* ribosomal gene query sequences, clearly showed an
2 18S, 5.8S and 28S rRNA gene block arrangement in okra. Two clusters are located at scaffolds ends,
3 though not coinciding with, or flanking telomere blocks. Another two clusters are positioned toward the
4 scaffold centre, and three scaffolds almost entirely consist of 18S-5.8S-28S gene clusters. These
5 scaffolds do not contain 5S rRNA clusters. Instead, 5S rRNA genes are organized in clusters separated
6 from the 18S-5.8S-28S units, clearly indicating an S-type rRNA gene arrangement (Goffová and Fajkus,
7 2021). In total we found four 5S rRNA clusters on four different scaffolds of which the largest cluster
8 consisted of almost 8,700 copies tandemly arranged on a single scaffold (Table S10). Signatures of
9 underlying chromosome evolution involving telomere fusion at ribosomal gene clusters were not
10 apparent though, as we did not encounter interstitial telomere repeats in rRNA clusters. *Genetic diversity*
11 *in okra*

12 To assess the genetic diversity in public okra germplasm, we used the okra reference genome to
13 call single nucleotide polymorphisms (SNPs) from several publicly available Illumina RNA-seq datasets
14 that we retrieved from the short-read archive (SRA) (<https://www.ncbi.nlm.nih.gov/sra>). Most of these
15 samples represent accessions originating from the Indian and Chinese parts of the Asian continent. The
16 average map rate for a panel of 11 samples was $93.2\% \pm 5.6\%$ (Data S1). The combined samples cover
17 20-25% of the reference genome, which is in line with the 26.6% genic portion of the genome and
18 suggesting a faithful structural annotation of the reference genome. The unusual high coverage for the
19 'Xianzhi' dataset (~65%) possibly was due to a deviating library preparation or DNA contamination, while
20 the lower coverage for the 'IIHR-299', 'Mahyco Arka Abhay', and 'Commercial' samples (10-12%) was
21 likely due to their smaller data size. The total panel size comprised 412,185 loci, for which 690,145 SNPs
22 were detected from coding regions of okra genes. The "Commercial" sample yielded only 1,741 unique
23 SNPs with reads mapping to 12.5% coverage of the reference genome. This SNP rate is substantially
24 lower compared to "Mahyco Arka Abhay" and "IIHR-299", while these 3 samples cover approximately
25 equal portions of the genome, suggesting that the 'Commercial' breeding line apparently shares a large
26 part of its ancestry with the reference cultivar. However, difference in tissue types, growth conditions,
27 data generation and processing workflows complicate direct sample comparison. Nevertheless, most of
28 the samples attain 20,000 unique SNPs and only 'Arka Anamika' and 'Danzhi' exceed this level. Although
29 we could not yet assess the genetic diversity in the non-genic portion, the okra accessions in the panel
30 apparently represent a low genetic diversity.

31 *Allopolyploid composition of okra*

32 We attempted to divide okra subgenomes by searching for patterns based on hierarchical
33 clustering of repeat k-mer counts without a supposition of ancestral species. We compared the clustering
34 results for okra with two k-mer test sets, of which the first comprised an artificially constructed hybrid

1 genome, consisting of merged tomato (*S. lycopersicum* cv. Heinz 1706) and a diploid potato (*S.*
2 *tuberosum* cv. Solyntus) genomes. The second set was generated from the allotetraploid cotton genome
3 (*Gossypium hirsutum*). As expected the cluster map for the artificial hybrid dataset separated into two
4 distinct subclusters, each representing repetitive k-mers from 12 tomato and potato chromosomes
5 (Figure S5). In addition, the cotton genome also clustered into two groups, clearly separating the
6 repetitive kmers generated from the subgenome A and D chromosomes (Hu *et al.*, 2019). Applying the
7 same method to the okra reference genome, of which five smallest scaffolds were removed, yielded two
8 distinct clusters of 30 and 43 scaffolds with a length of 636 Mbp and 557 Mbp, respectively (Figure S6).
9 Scaffolds from cluster 1 are overall larger than from cluster 2 (Figure S7). Apparently the repeat k-mer
10 pattern for okra points to two distinct subgenomes and, together with the apparent absence of
11 multivalent pairing of metaphase chromosomes, suggests an allotetraploid nature of *Abelmoschus*
12 *esculentus*. In addition, we could not find clear evidence of erosion, as the clusters had comparable
13 BUSCO completeness scores of 90.2% and 89.3%, with duplicate rates of 21.8% and 19.4%
14 respectively. This suggests a relatively recent hybridization of ancestral species, yet without clear
15 evidence for an emerging dominant subgenome.

16 We subsequently aligned the clustered scaffolds to further investigate the orthology between the
17 two subgenomes. Although we found only partial alignments and several inversions, there is substantial
18 homoeology between scaffolds (figure S8), confirming in general there are no homoeologs within a single
19 cluster. Specifically, the network of scaffolds displayed the portion of BUSCO genes shared between
20 scaffolds (figure S9), indicates the strong homoeology between scaffolds of cluster 1 and 2,
21 corroborating the alignment dot-plot.

22 Further repeat annotation reveals only 784 out of 22,100 repeat classes are present in >90% of
23 scaffolds in both clusters. A small number of 14 repeat classes are present in more than 90% of the
24 scaffolds in either cluster, while occurring in less than 10% of scaffolds in the other. These include cluster
25 specific repeat classes, occurring in high copy numbers in cluster 1 but not in cluster 2 and vice versa
26 (Figure S10). Although such repeats do not show similarity to known repeats and we can only speculate
27 about their origin, they are clearly related to different ancestral progenitor species, further pointing to an
28 allotetraploid nature of okra.

29 *Candidate genes assigned to phenylpropanoid, flavonoid, and flavone and flavonol biosynthesis pathways*

30 Polyphenols represent one of the most ubiquitous class of secondary metabolites in okra fruits.
31 An important subclass of polyphenols are flavonoids, of which myricetin, quercetin, isoquercitrin and
32 quercetin-3-O-gentiobioside derivatives have been implicated in antidiabetic activity (Liu *et al.*, 2005; Lei

1 *et al.*, 2017; Wu *et al.*, 2020; Peter *et al.*, 2021). Myricitin was previously detected in *Abelmoschus*
2 *moschatus* (Liu *et al.*, 2005). Recently, the bioactive phytochemicals isoquercitrin and quercetin-3-O-
3 gentiobioside, and to a lesser extent also rutin and catechin, were detected as the major phenolic
4 compounds in okra fruits (Wu *et al.*, 2020). Their biosynthesis in the flavonoid and, flavone and flavonol
5 biosynthesis pathways (KEGG reference pathways 00941 and 00944) is thought to start with p-
6 coumaroyl-CoA and cinnamoyl-CoA precursors that are synthesized in the phenylpropanoid pathway
7 (KEGG reference pathway 00940). To find putative enzyme coding okra genes that may function in
8 phenylpropanoid, flavonoid, flavone and flavonol biosynthesis, 142,571 extracted amino acid query
9 sequences from predicted okra genes and putative splice variants were mapped against the manually
10 curated KEGG GENES database, using the KEGG Automatic Annotation Server (KAAS) ([KAAS - KEGG](#)
11 [Automatic Annotation Server \(genome.jp\)](#) (Moriya *et al.*, 2007). We identified 33,641 amino acid
12 sequences that could be assigned to 395 KEGG metabolic pathway maps based on a best bi-directional
13 hit (BBH). Currently, in total there are N=1302 manually annotated from *Malvaceae* species *Theobroma*
14 *cacao* (cacao), *Gossypium arboreum*, *Gossypium hirsutum* (cotton), *Gossypium raimondii*, and *Durio*
15 *zebithinus* (durian) of which K=99 enzyme coding reference genes are known for the phenylpropanoid
16 (K₁=36), flavonoid (K₂=30), or flavone and flavonol (K₃=33) biosynthesis pathway in KEGG. Of the
17 33,641 okra amino acid queries, n=47 putative okra orthologs were assigned to the KEGG
18 phenylpropanoid biosynthesis (n₁=16), flavonoid biosynthesis (n₂=17) and flavone and flavonol
19 biosynthesis (n₃=14) pathways respectively, adding up to n=41 distinct putative okra enzyme orthologs.
20 We subsequently assessed the mapping probability of okra orthologs to the reference pathways based on
21 the known *Malvaceae* enzymes and the okra BBH. Mapping confidence values p₁=1.43e⁻¹², p₂=1.15e⁻¹²,
22 and p₃=0.0 pointed to a confident assignment of okra orthologs to phenylpropanoid biosynthesis,
23 flavonoid biosynthesis, and flavone and flavonol biosynthesis pathways respectively. Copy numbers for
24 putative genes possibly involved in the conversion p-coumaroyl-CoA and cinnamoyl-CoA precursors
25 varied extensively. Only a single putative gene orthologous to a 5-O-(4-coumaroyl)-D-quinic_3'-
26 monooxygenase (EC:1.14.14.96) from *Durio zebithinus* (XP_022742205) with an amino acid identity of
27 93.9% was found, whereas 14 putative orthologs to shikimate O-hydroxycinnamoyltransferase
28 (EC2.3.1.133) were detected, with a highest amino acid identity (94.9%) to the ortholog from
29 *Gossypium arboreum* (XP_017607223). The coverage of okra orthologs mapped to these biosynthesis
30 pathways is shown in figures 3 and S11. The alternative metabolic routes, leading to the biosynthesis of
31 quercetin, myricitin, isoquercitrin, and quercetin-3-O-gentiobioside derivatives in these pathways, involve
32 three critical flavonol synthases CYP74A (flavonoid 3',5'-hydroxylase, EC:1.14.14.81), CYP75B1 flavonoid
33 3'-monooxygenase (EC:1.14.14.82) and (FLS) dihydroflavonol,2-oxoglutarate:oxygen oxidoreductase
34 (EC:1.14.20.6). Apparently, putative orthologs for CYP74A and CYP75B1 are encoded by a single copy

1 gene in okra, whereas the FLS oxidoreductase catalytic activity, that is thought to catalyse the
2 conversion of several dihydroflavonol intermediates into quercitin and myricitin, appears represented by
3 5 putative okra homologs, suggesting that the conversion into quercitin, isoquercitrin, rutin and catechin
4 mainly runs via a dihydrokaempferol intermediate.

5 **Conclusions**

6 We successfully applied multiple DNA sequencing and scaffolding techniques to reconstruct the
7 65 chromosomes of a haploid okra sibling. With a final result of 49 scaffolds with telomeres at both ends,
8 illustrating chromosome completeness, and 16 scaffolds with at least 1 telomere end, representing
9 partially resolved chromosomes, the assembly contiguity for okra surpasses that of other crop genome
10 assemblies. PacBio high fidelity (HiFi) reads and availability of a haploid sample were key to the high
11 quality of the reconstructed okra genome. The more pronounced diploid behavior of okra, exemplified by
12 the smudgeplot spectrum showing 75% AB dominance, when compared to the k-mer spectrum for
13 allotetraploid cotton with 50% AB dominance, illustrates a decreased complexity that benefitted the
14 genome reconstruction. This diploid-like nature is also in line with our cytogenetic observations on pollen
15 mother cells at pachytene, showing chromosomes strikingly diploid-like with clear bivalents. The okra
16 genome is characterized by a high chromosome count and relatively small chromosomes that do not
17 extend beyond 30 Mbp in length. The total haploid assembly size of ~1.2 Gbp apparently is consistent
18 with the slightly larger k-mer based haploid genome size estimation. Mapping of telomeric sequences and
19 gene dense sections toward both scaffold ends, together with repeat dense sections mapping toward
20 more centrally located scaffold sections, is in line with the observed chromatin landscape at pachytene
21 for several okra chromosomes that display pericentromeric heterochromatin and distal euchromatin,
22 which have been shown repeat-rich and gene-rich respectively also in many species genomes. The
23 overall repeat content (57%) is lower than described for *Gossypium raimondii* (70.7%), but higher than
24 for *Theobroma cacao* (29.4%) (Novák *et al.*, 2020), whereas these genomes are smaller. Our repeat
25 classification remains incomplete though, due to a lack of diversity in annotated repeat libraries.
26 Furthermore, with over 130,000 putative genes, the gene prediction count seems inflated, compared to
27 other species genomes like *Arabidopsis* and rice (*Oryza sativa*), containing around 38,000 and 35,000
28 coding, non-coding and pseudogenes respectively. Some increase may be attributed to
29 allopolyploidization. On the other hand, a decreased selective pressure may cause many genes to
30 accumulate *de novo* mutations and convert into pseudogenes (Bird *et al.*, 2018). Nevertheless, over 88%
31 of the predicted exons was supported by high-quality long read data and predicted proteins mostly
32 returned partial matches to several different databases, substantiating our gene prediction. In this
33 respect, the structural and functional annotation revealed putative enzyme coding genes that we could

1 map to phenylpropanoid, flavonoid, flavone and flavonol metabolic pathways, likely underlying the
2 biosynthesis of an array of secondary metabolites that have been implicated in dietary and therapeutic
3 bioactivity.

4 Identification of subgenomes from distinct k-mer repeat profiles point to an allotetraploid
5 composition of the genome. The two subgenomes apparently have unequal chromosome counts, though
6 are similar in total sequence length. The BUSCO scores suggest nearly complete subgenomes, likely
7 reflecting a relatively recent hybridization of two progenitor species without clear evidence for
8 subgenome dominance. Substantial numbers of duplicated BUSCO genes in the subgenomes might
9 reflect hybridization events prior to allotetraploidization. In this respect the additional grouping of k-mers
10 within each of the two k-mer repeat clusters is suggestive and intriguing, though need more detailed
11 analysis. Finally, we conclude that the annotated high-quality genome provides a solid basis for further
12 okra breeding application, including diversity and compatibility screening, and marker development.

13

14

15

16

17

18

19

20

21

22

23

24

25

26

1 **Acknowledgements**

2 We wish to thank Hortigenetics Research of East West Seed (S.E. Asia) Ltd., ENZA Zaden Research and
3 Development B.V., Genetwister Technologies B.V., Nunhems Netherlands B.V., Syngenta Seeds B.V.,
4 Takii & Company Ltd., HM.Clause, SA., UPL Ltd., Namdhari Seeds Pvt. Ltd., Maharashtra Hybrid Seeds
5 Co. Pvt. Ltd. and Acsen HyVeg Pvt. Ltd. for providing material and support to the okra genome project.

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

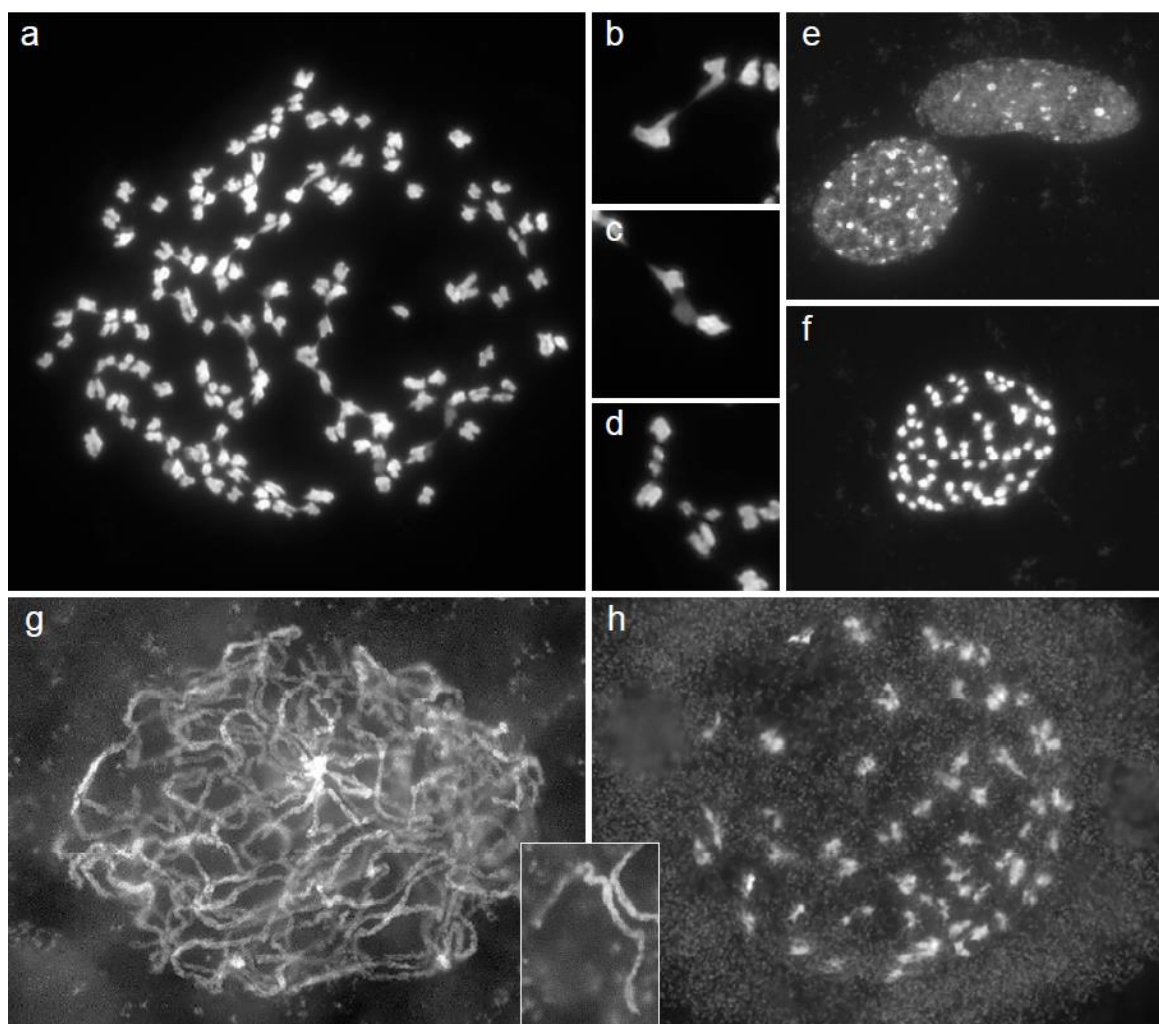
21

22

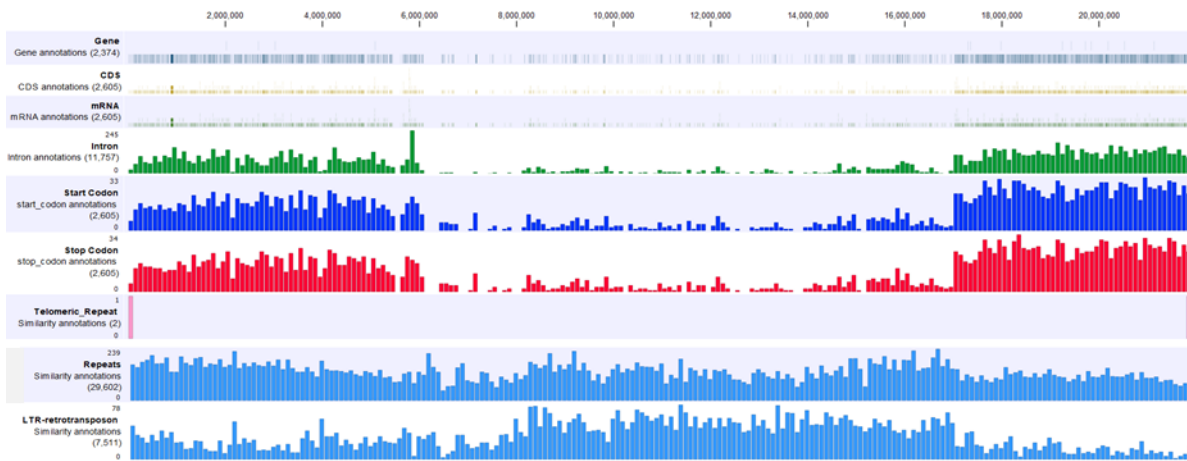
23

24

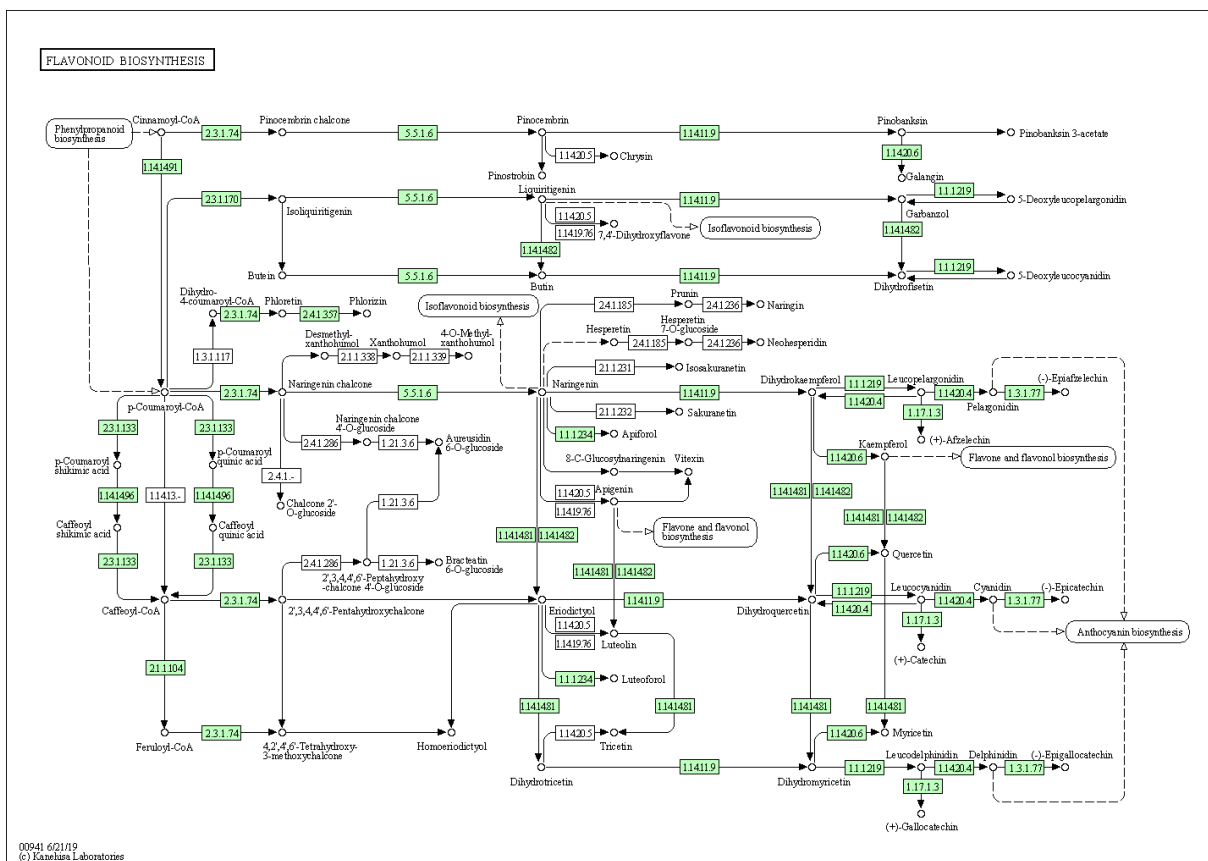
1 **Figures**



3 **Figure 1:** Mitotic cells in root tip meristems. (a) Example of a well spread metaphase complement of a
4 diploid okra plant with $2n=130$. (b) Magnification of two chromosomes that are joined by telomere
5 connectives. (c) Chromosome pair with a less fluorescing region, likely representing a decondensed
6 Nucleolar Organizer Region (NOR). (d) chromosomes with small satellites. (e) Two interphase nuclei
7 displaying a striking difference in number of highly condensed chromocenters, regions of the
8 pericentromere heterochromatin and NORs. The top nucleus has about 10 chromocenters; some other
9 nuclei can have more than 100 of such condensed regions. (f) Metaphase complement of a haploid okra
10 plant ($2n=65$). (g,h) Meiotic chromosomes in pollen mother cells of a diploid okra plant. (g) Cell at
11 pachytene stage. Most of the chromosomes are fully and regularly paired without clear indications for
12 multiple synapsis, pairing loops or pairing partner switches. The brightly fluorescing regions are the
13 pericentromeres, see also the inset between the figures g and h. (h) Cell at diakinesis. A greater part of
14 the chromosomes clearly forms bivalents. Magnification bars in the figures equals $5 \mu\text{m}$.



1
2 **Figure 2:** Structural annotation. Annotation feature classes for a 25 Mbp okra scaffold and coordinate
3 positions are indicated at the left and top side of the plot respectively. Bar heights in each row
4 corresponds to the relative frequency of genic, non-genic, and repeat class per scaffold segment of
5 approximately 115 kb.



6
7 **Figure 3:** The flavonoid KEGG bio-synthesis pathway in *Abelmoschus esculentus*. Putative okra enzyme
8 coding genes for which a bi-directional best hit was found to reference pathway enzymes are shown with
9 coloured EC identifiers.

10

1 **Tables**

| Assembly statistic | Primary ctgs | Alternative ctgs | Hybrid scfds |
|---------------------------|---------------------|-------------------------|---------------------|
| Ctgs/scfds | 3,051 | 972 | 80 (78) |
| Total length | 1,223.6 Mb | 31.0 Mb | 1,194.5 Mb |
| Median length | 32.4 kb | 24.5 kb | 16.320 Mb |
| Max length | 25.1 Gb | 652 kb | 29.444 Mb |
| Min length | 13.3 kb | 10.3 kb | 125 kb |
| N50 length | 10.6 Gb | 30.5 kb | 18.929 Mb |
| N50 index | 43 | 126 | 27 |
| N95 length | 483 kb | 15.5 kb | 7.206 Mb |
| N95 index | 168 | 465 | 64 |
| GC content | 34.36% | 47.6% | 33.76% |

2

3 **Table 1:** NGS assembly and hybrid scaffolding statistics. Sequences were assembled using the Hifiasm
 4 assembler and scaffold with Bionano Genomics genome maps. Number of scaffolds obtained with 10X
 5 Genomics linked reads is indicated between brackets.

| Class | BUSCO statistics | | | | |
|-------------------------|-------------------------|---------------|---------------|---------------|---------------|
| Assembly version | 1 | 2 | 3 | 4 | 4 |
| Coverage | 20X | 84X | 95X | 95X | 95X |
| Ctgs/scfds | all | all | all | primary | alternative |
| Complete | 2,288 (98.3%) | 2,271 (97.7%) | 2,269 (97.6%) | 2,288 (98.4%) | 66 (2.8%) |
| Single copy | 266 (11.4%) | 311 (13.4%) | 313 (13.5%) | 284 (12.2%) | 66 (2.8%) |
| Multiplicated | 2,022 (86.9%) | 1,960 (84.3%) | 1,956 (84.1%) | 2,004 (86.2%) | 0 (0%) |
| Duplicated | n.d | n.d | n.d | 1,150 (49.4%) | 0 (0%) |
| Triplicated | n.d | n.d | n.d | 843 (36.2%) | 0 (0%) |
| Quadruplicated | n.d | n.d | n.d | 7 (0.3%) | 0 (0%) |
| Quintuplicated | n.d | n.d | n.d | 3 (0.1%) | 0 (0%) |
| Sextuplicated | n.d | n.d | n.d | 0 (0%) | 0 (0%) |
| Septuplicated | n.d | n.d | n.d | 1 (0.04%) | 0 (0%) |
| Fragmented | 5 (0.2%) | 5 (0.2%) | 6 (0.3%) | 6 (0.3%) | 9 (0.4%) |
| Missing | 33 (1.5%) | 50 (2.1%) | 51 (2.1%) | 32 (1.3%) | 2,251 (96.8%) |

| | | | | | |
|--------------|-------|-------|-------|-------|-------|
| Total | 2,326 | 2,326 | 2,326 | 2,326 | 2,326 |
|--------------|-------|-------|-------|-------|-------|

1

2 **Table 2:** Detection of ortholog core genes. Genome assemblies at different coverage levels were
3 analysed to assess the assembly completeness. BUSCO classes are shown as single copy or multiplied
4 ortholog. Multiplied orthologs are subdivided into additional copy classes as indicated. For each
5 assembly coverage level BUSCO counts in primary, alternative, and all contigs are shown in absolute
6 numbers and percentages of total expected orthologs (between brackets), or n.d (not determined).

| Class | Count | Av. size | Total length |
|-----------------|--------------|-----------------|------------------------|
| Total Scfds | 4,023 | 328,851 | 1,322,968,356 (100%) |
| Large Scfds | 78 | 14,932,447 | 1,194,595,770 (90.30%) |
| Gene | 130,324 | 2,537 | 330,639,435 (24.99%) |
| CDS | 150,032 | 2,497 | 374,629,904 (28.32%) |
| mRNA | 150,032 | 2,497 | 374,629,904 (28.32%) |
| Start | 150,004 | 3 | - |
| Stop | 150,009 | 3 | - |
| Intron | 676,681 | 307 | 207,741,067 (15.70%) |
| sRNA | 2308 | 797 | 1,839,960 (0.139%) |
| Total repeats | 1,351,943 | - | 677,354,628 (51.20%) |
| unclassified | 834,282 | 321 | 268,086,453 (20.26%) |
| [TTTAGGG]n | 123 | 1,810 | 222,579 (0.017%) |
| Retroelements | 442,480 | 858 | 379,556,626 (28.69%) |
| LTRs | 389,339 | 924 | 359,876,901 (27.20%) |
| Gypsy | 146,673 | 1,376 | 201,862,257 (15.26%) |
| Ty1/Copia | 122,317 | 975 | 119,289,622 (9.02%) |
| LINES | 26,571 | 650 | 17,281,993 (1.31%) |
| SINES | 312 | 507 | 158,078 (0.01%) |
| DNA transposon | 46,849 | 478 | 15,456,661 (1.17%) |
| Hobo-Ac | 16,781 | 354 | 5,941,881 (0.45%) |
| Tc1-Pogo | 645 | 212 | 136,523 (0.01%) |
| 5S rDNA | 9644 | 90 | 871,856 (0.066%) |
| 5S rDNA partial | 129 | 31 | 4,035 (<0.001%) |
| 5.8S rDNA | 201 | 622 | 125,073 (0.009%) |

| | | | |
|-------------------|-----|-------|------------------|
| 5.8S rDNA partial | 23 | 231 | 5,311 (<0.001%) |
| 18S rDNA | 183 | 1,750 | 320,241 (0.024%) |
| 18S rDNA partial | 170 | 372 | 63,169 (0.005%) |
| 28S rDNA | 167 | 3,368 | 562,376 (0.043%) |
| 28S rDNA partial | 266 | 628 | 167,176 (0.013%) |

1

2

Table 3: Structural annotation for the 78 largest okra scaffolds. Features are classified into genic and

3

repeat elements as indicated. Statistics are in nucleotide length and in fractions of total scaffold length.

4

5

6

7

8

9

10

11

12

13

14

15

16

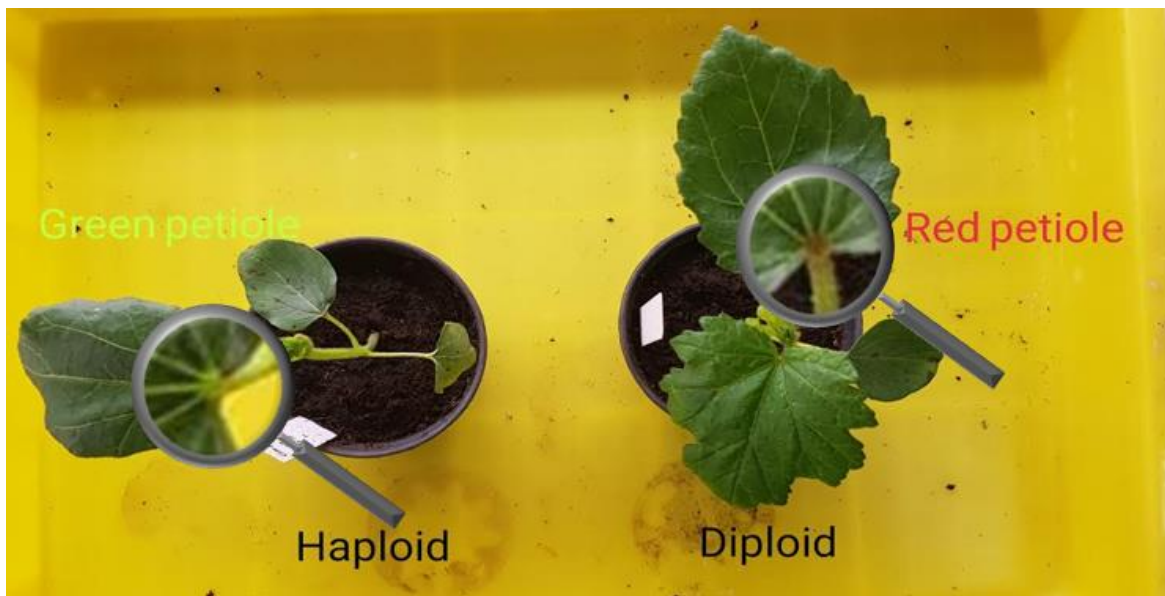
17

18

19

20

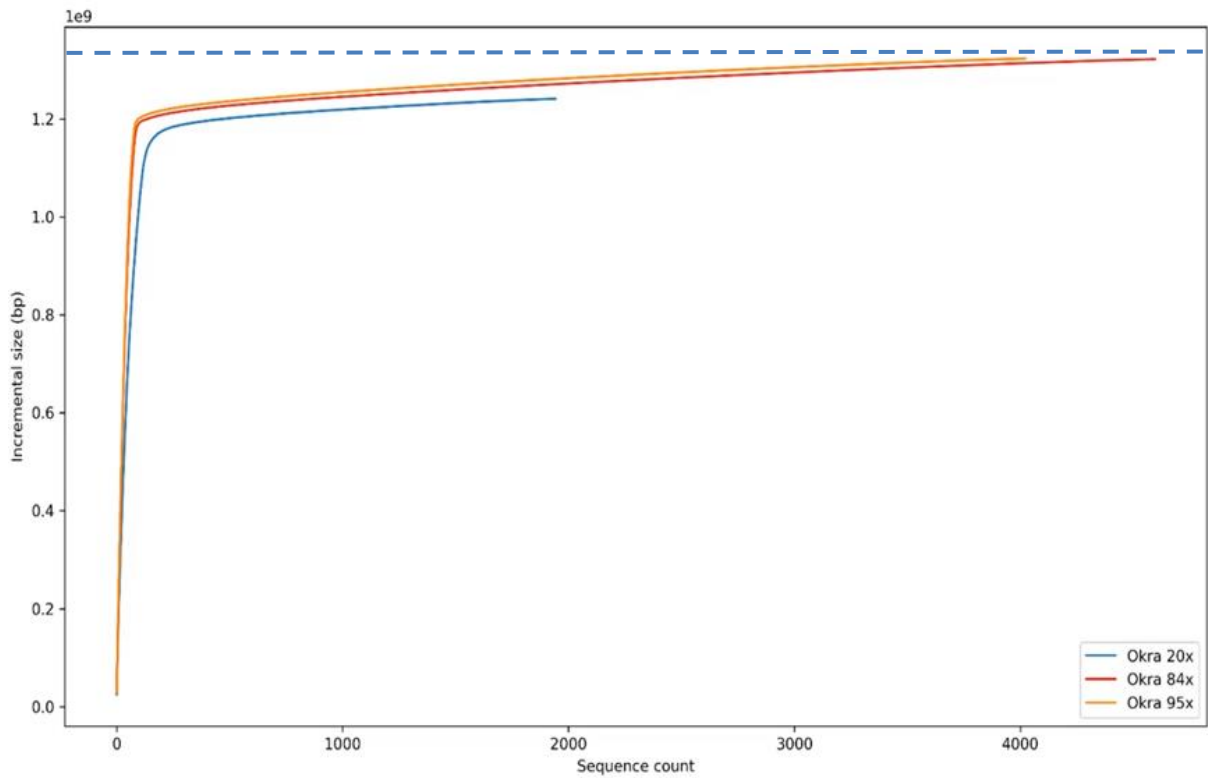
1 **Supplementary figures**



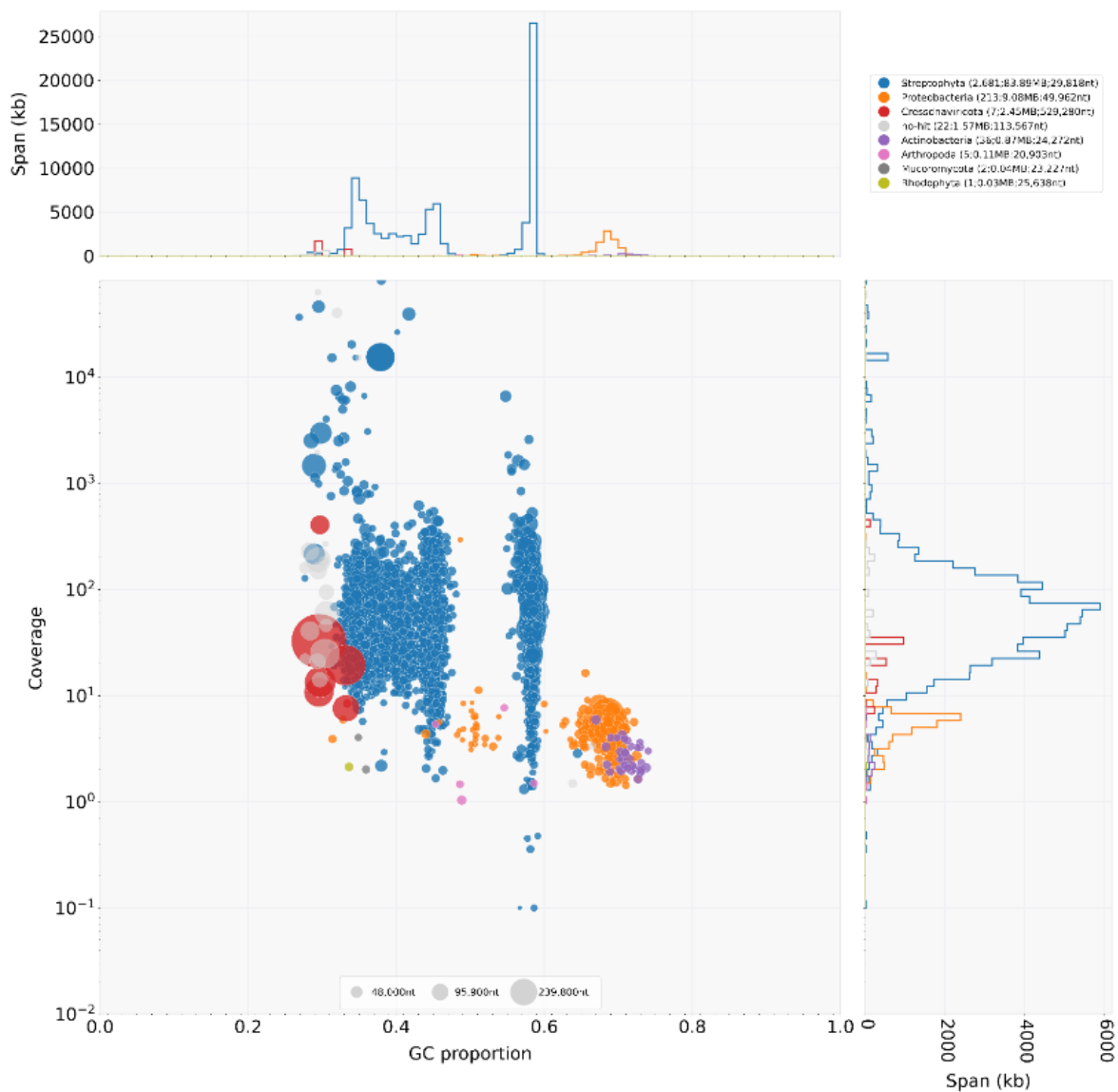
3 **Figure S1:** Phenotypes of diploid and haploid Okra plants. The magnifying glass in the image is placed
4 over the position of the green and the red petiole for the haploid (left) and diploid (right) Okra plant
5 respectively.

6

7



1
2 **Figure S2:** The incremental genome assembly size for Okra. The A50 plot for contigs larger than 100 bp
3 shows the assembly size in Gbp on the y-axis is plotted against the incremental contig count at 20X, 84X
4 and 95X sequence coverage indicated by the light blue, red and orange curve respectively.



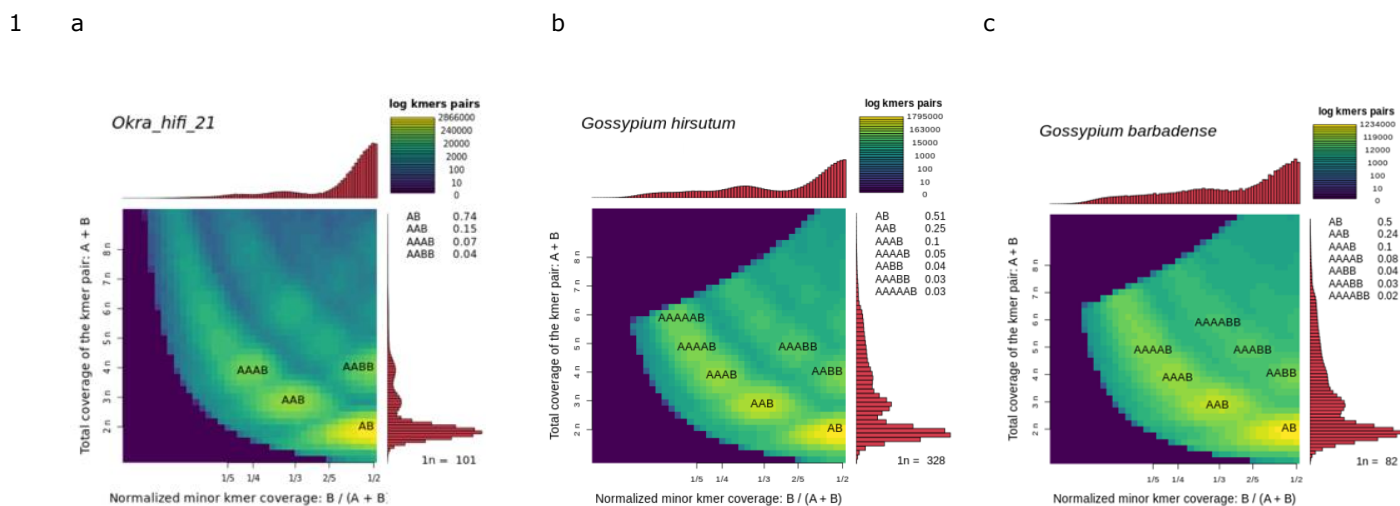
1
2 **Figure S3:** Taxon annotated GC coverage plot. In the low left panel the proportion of GC bases (x-axis)
3 and read coverage (y-axis) for 527 alternative contigs are shown. Each coloured dot in the graph
4 corresponds to a contig. Colours correspond to species classes for which a best BlastN match was found
5 in annotated databases. In the top left graph and the low right graph the relative proportion for each
6 class is depicted with respect to the coverage and GC content, for which colour codes match species
7 classes as indicated in the legend at the top right.

8

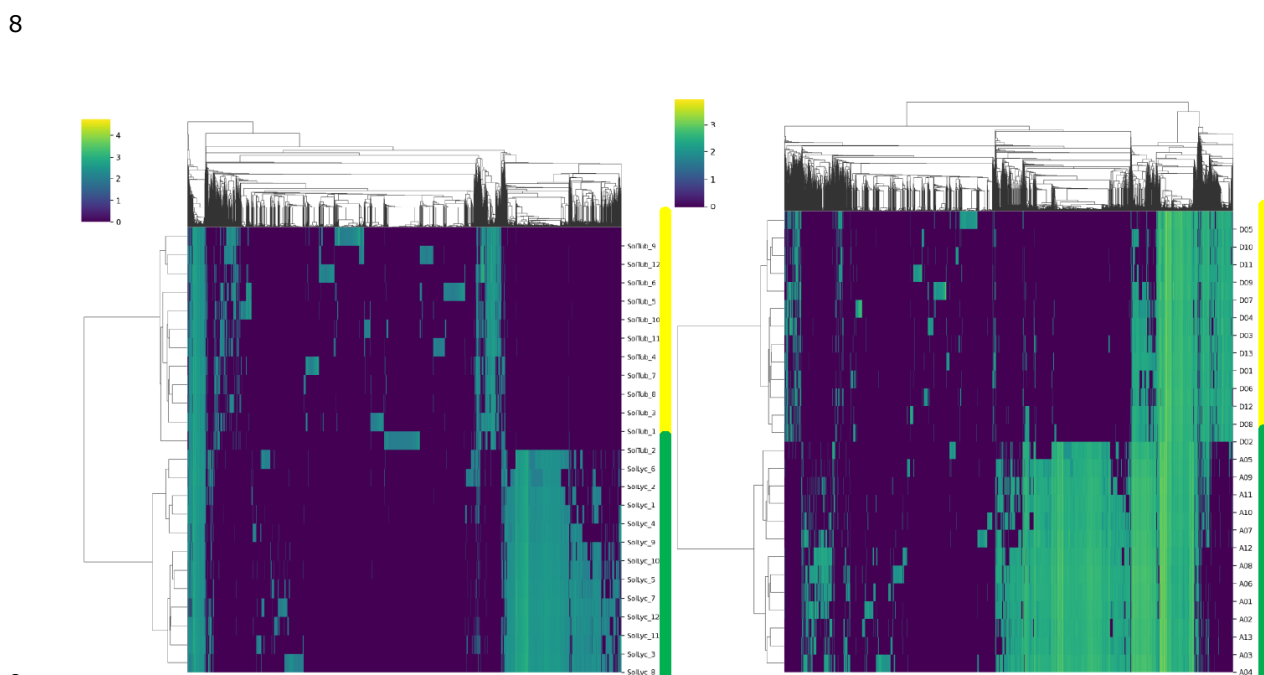
9

10

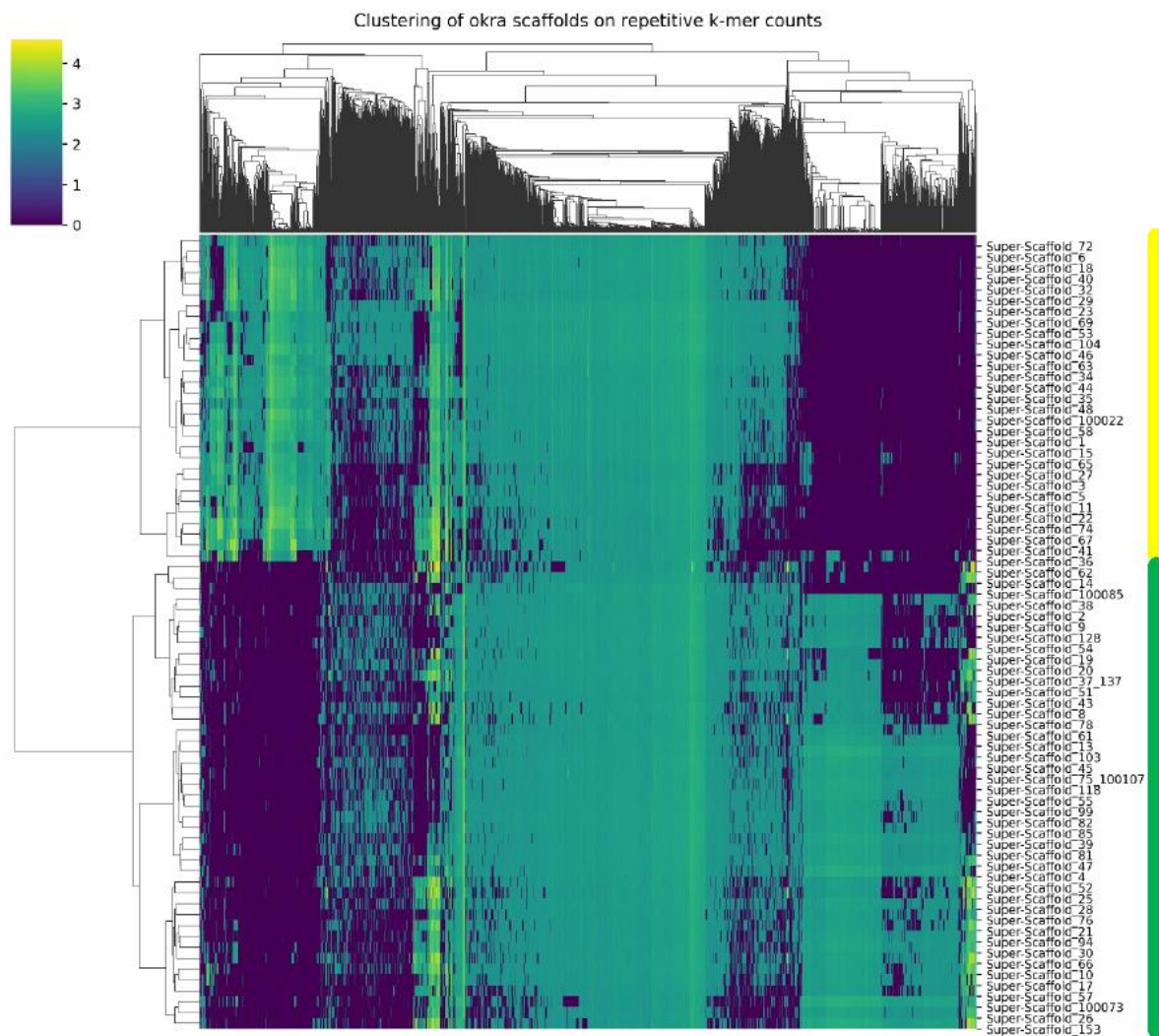
11



3
4 **Figure S4:** Smudgeplots for and haploid Okra (*Abelmoschus esculentus*) (a), and two allotetraploid
5 cotton species *Gossypium hirsutum* (b), and *Gossypium barbadense* (c). Smudgeplots are shown in log
6 scale. The coloration indicating the approximate number of k-mer pairs per bin and the fraction of each
7 kmer type is indicated in the top right legend of each plot.



9
10 **Figure S5:** Cluster maps of repetitive 13-mer counts. K-mers generated for an artificial hybrid genome
11 constructed from merged genomes of *S. lycopersicum* SL.4.0 and *S. tuberosum* cv. Solyntus (left panel)
12 and from the allotetraploid cotton *G. hirsutum* genome. Yellow and green bars next to the chromosome
13 identifiers mark the potato and tomato chromosomes (left panel) and the cotton chromosomes from the
14 A and D subgenomes (right panel) respectively. Color code bar indicates log10 scaled kmer counts.



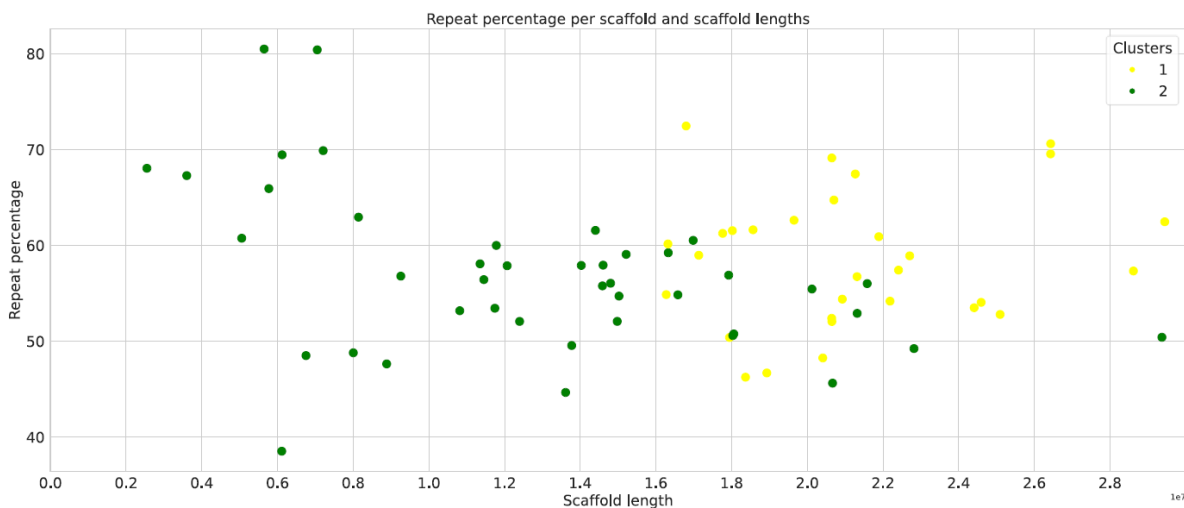
1

2 **Figure S6:** Cluster maps of repetitive 13-mer counts for the okra reference genome. Yellow and green

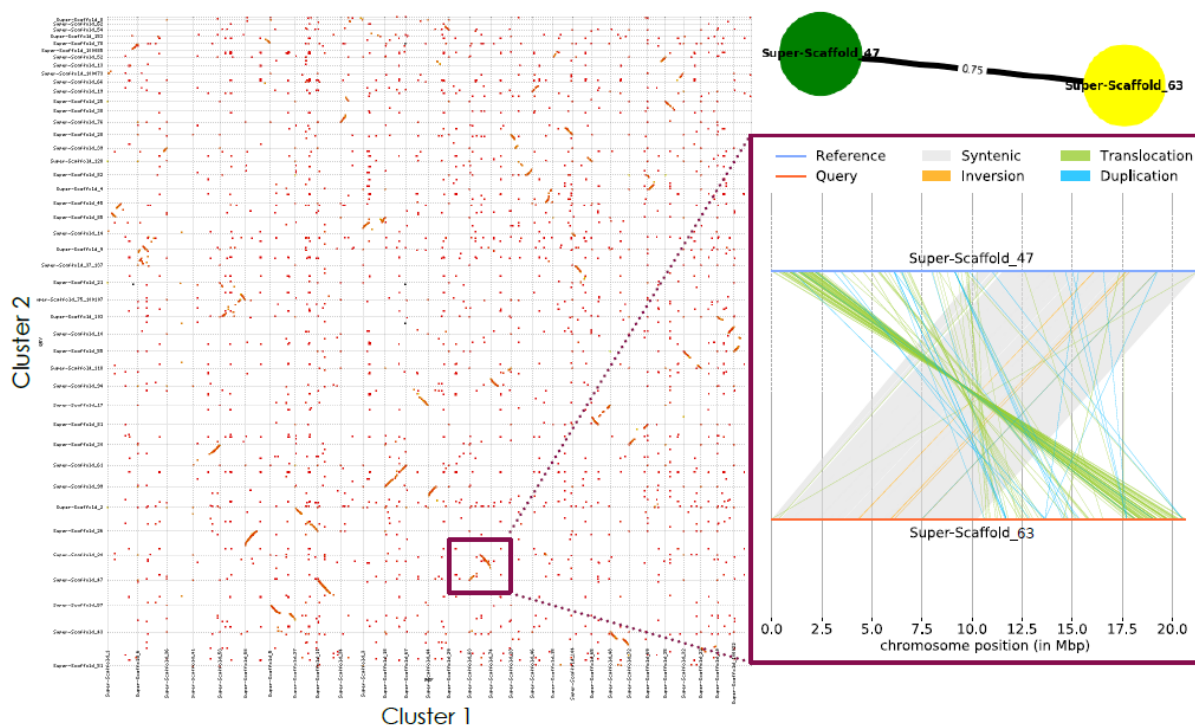
3 bars next to the identifiers mark the superscaffolds clustering in cluster 1 and 2 respectively. Color code

4 bar indicates log₁₀ scaled kmer counts. A clear separation between cluster 1 and 2 based on repeat

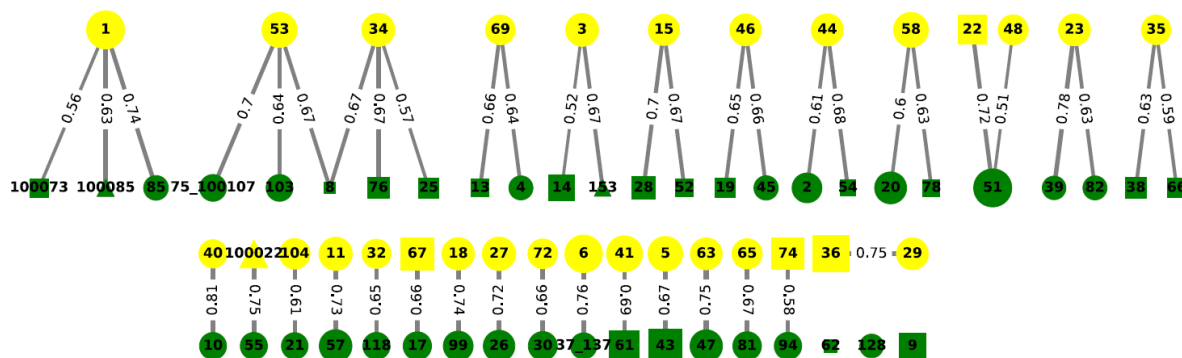
5 count and distinct repeat profile is apparent.



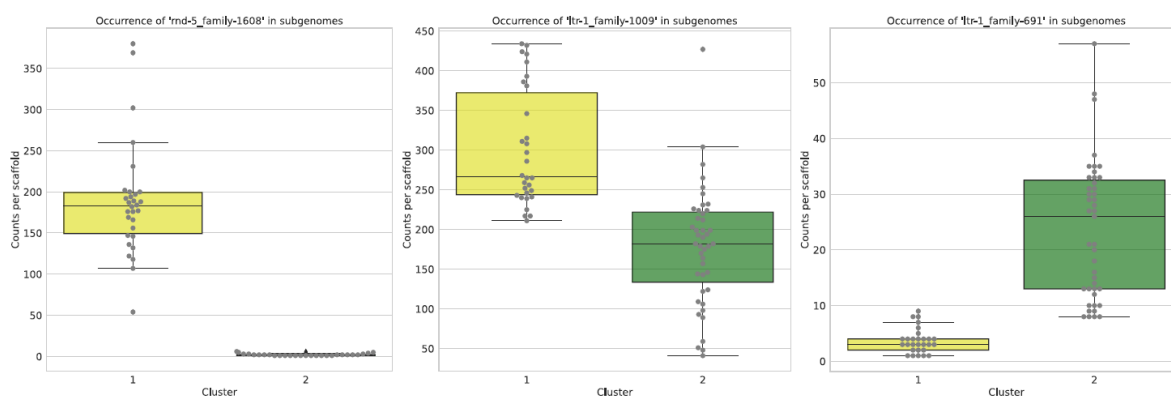
1
 2 **Figure S7:** Repeat content analysis in superscaffolds. Superscaffolds assigned to cluster1 and 2 are
 3 represented by yellow and green dots respectively, and have been separated by scaffold length (x-axis)
 4 and repeat percentage (y-axis).



5
 6 **Figure S8:** Dot plot alignment of superscaffolds from subgenomes. Superscaffolds are assigned to
 7 cluster 1 or 2 according to their kmer clustering profile. The top right graph shows two homoeologous
 8 superscaffolds 63 (cluster 1) and 47 (cluster 2) having 75% of BUSCO genes in common. The bottom
 9 right alignment detail of the aforementioned superscaffolds are partially syntenic, sharing a large
 10 inversion.



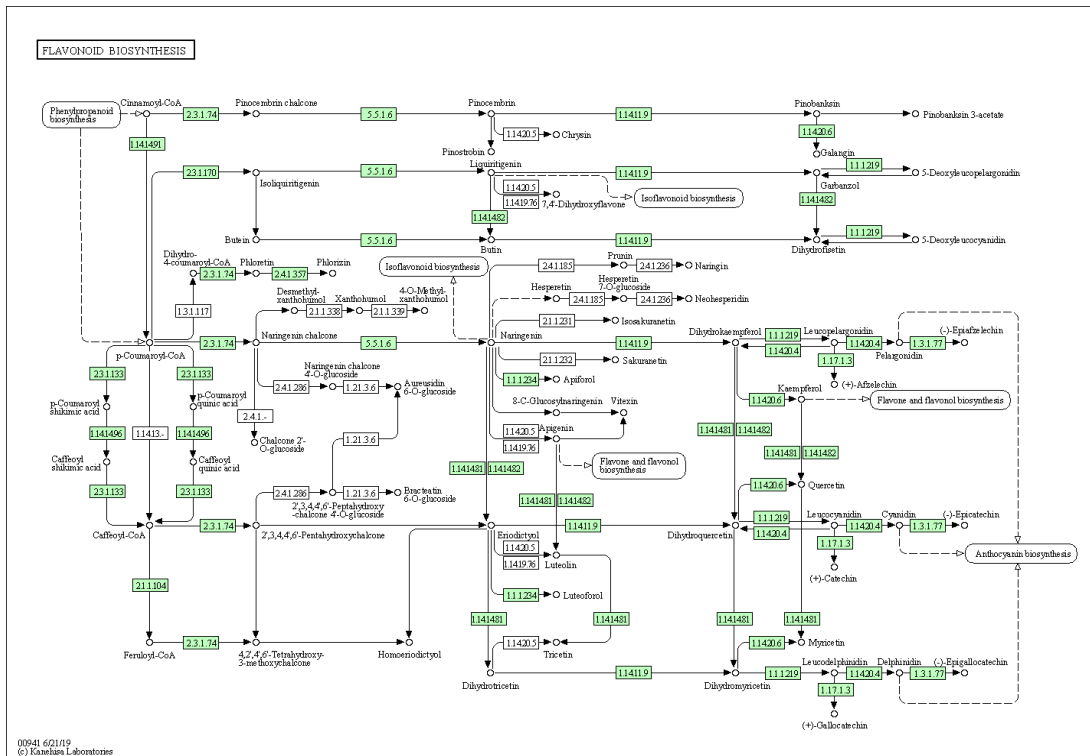
1
2 **Figure S9:** BUSCO connectivity graph. Yellow and green color coded nodes correspond to scaffolds from
3 cluster 1 and 2 respectively. Edges between the nodes indicate the percentage of shared BUSCO genes
4 between each scaffold pair. Note that a single node can have multiple edges. Pairs of scaffolds point to
5 links of homoeology between scaffolds.



7
8 **Figure S10:** Repeat count of cluster specific and shared repeats between subgenomes. Occurrence of 3
9 distinct repeat families are shown as counts per scaffold (y-axis) that are divided over distinct clusters 1
10 and 2 (x-axis). Counts per scaffold are represented by grey dots. The repeat family identifier is indicated
11 above each plot. The left panel shows the occurrence of an unclassified repeat family in cluster 1 specific
12 scaffolds while absent in cluster 2 scaffolds. The right panel shows an unclassified cluster 2 specific
13 repeat family. The unclassified repeat family in the middle graph is cluster unspecific.

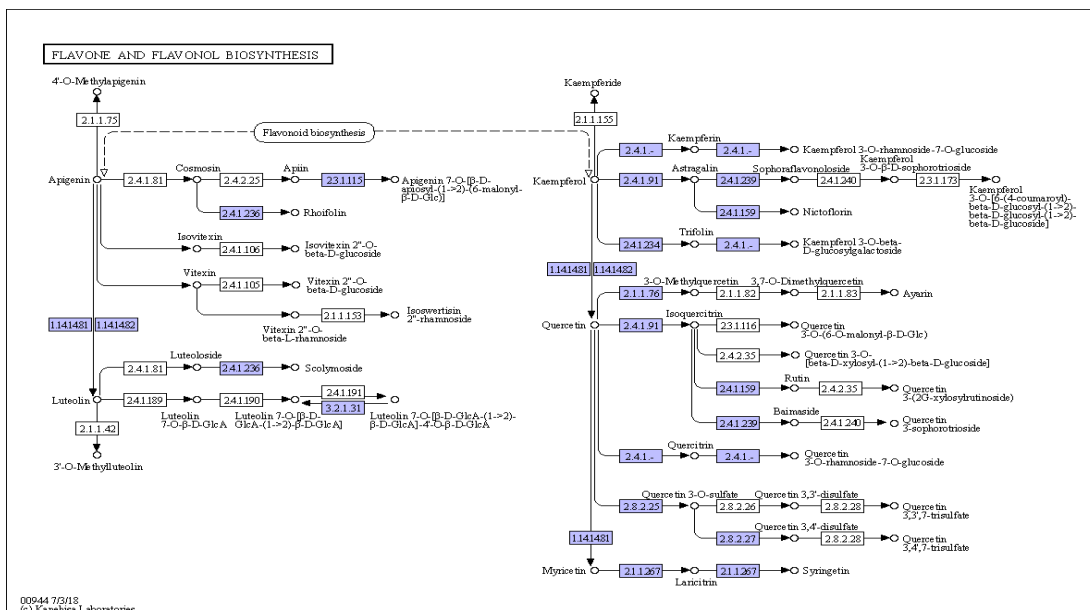
14
15
16
17

1 **A**



2

3 **B**



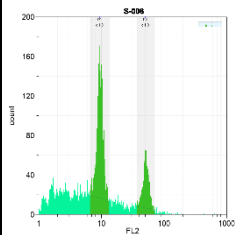
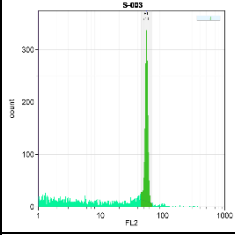
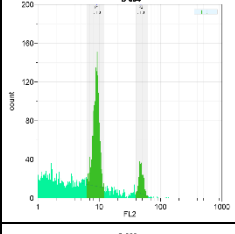
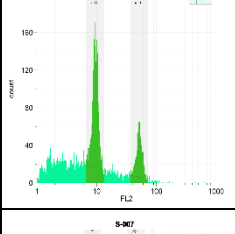
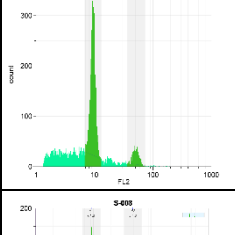
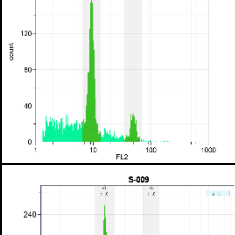
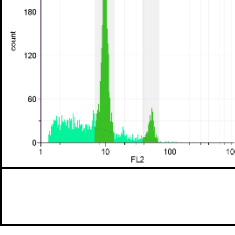
4

5 **Figure S11:** The flavonoid (A), and flavone and flavonol (B) KEGG bio-synthesis pathways in
 6 *Abelmoschus esculentus* (<http://www.kegg.jp/kegg/kegg1.html>). Putative okra enzyme coding genes for
 7 which a bi-directional best hit was found to enzymes pathways are shown with coloured EC identifiers.

8

9

1 **Supplementary tables**

| Sample ID | Info | Rel. DNA amount | Flow Histogram |
|----------------|-------------|-----------------|--|
| S-002 | Okra1 | |  |
| S-003 | Agave | 15.90 |  |
| S-004 | Okra1+Agave | 2.99 |  |
| S-006 | Okra2+Agave | 2.94 |  |
| S-007 | Okra3+Agave | 2.94 |  |
| S-008 | Okra4+Agave | 3.02 |  |
| S-009 | Okra5+Agave | 3.05 |  |
| Average | | 2.99 (±0.01) | |

1 **Table S1:** DNA amount of nuclei samples from okra root tip cells. DNA amount of okra replicate samples
 2 in picogram quantities was compared to a reference sample from *Agave Americana*. In the right column
 3 flow histograms of Okra samples and the Agave reference sample are shown. The count of observed
 4 nuclei in each histogram is depicted on the y-axis and is proportional to fluorescent intensity of each
 5 peak. The position of the peak along the x-axis is proportional to the relative DNA amount in each nuclei.
 6

| Library type | Read length (bp) | raw yield (Gbp) | Read 1 >Q30 (bp) | Read 2 >Q30 (bp) |
|--------------------|---------------------------------|-------------------------|----------------------|--------------------|
| 10X PE | 151 | 362 | 91.7 | 89.0 |
| 10X PE | 151 | 318 | 92.9 | 90.4 |
| 10X PE | 151 | 122 | 92.5 | 88.9 |
| Library type | Av. polymerase read length (bp) | Subread length N50 (bp) | Longest subread | |
| Pacbio HiFi | 411.39 | | 16,927 | |
| Pacbio HiFi | 516.25 | | 13,086 | |
| Pacbio HiFi | 472.40 | | 12,732 | |
| Library type | Total molecules | Total length (bp) | Average length (kbp) | Molecule N50 (kbp) |
| BioNano unfiltered | 77,407,575 | 4,887,375.86 | 63.138 | 90.175 |
| BioNano filtered | 5,500,210 | 1,184,887.6 | 215.426 | 206.88 |

7
 8 **Table S2:** Genome sequencing and genome map data statistics. Linked read sequencing for three 10X
 9 Genomics libraries was performed using Illumina paired-end (PE) sequencing. Circular consensus
 10 sequencing (CCS) was performed for three Pacbio Hifi sequence libraries. Genome map data was
 11 produced for one BioNano DLE labelled library.

12

13

14

| Rank | No. hits | Hits | Family |
|------|----------|---|----------------------|
| 1 | 319 | Eukaryota <i>Gossypium hirsutum</i> | <i>Malvaceae</i> |
| 2 | 224 | Eukaryota <i>Theobroma cacao</i> | <i>Malvaceae</i> |
| 3 | 80 | Eukaryota <i>Gossypium raimondii</i> | <i>Malvaceae</i> |
| 4 | 72 | Eukaryota <i>Gossypium arboreum</i> | <i>Malvaceae</i> |
| 5 | 62 | Eukaryota <i>Durio zibethinus</i> | <i>Malvaceae</i> |
| 6 | 57 | No hits found | - |
| 7 | 29 | Eukaryota <i>Abelmoschus esculentus</i> | <i>Malvaceae</i> |
| 8 | 20 | Eukaryota <i>Hibiscus cannabinus</i> | <i>Malvaceae</i> |
| 9 | 14 | Eukaryota <i>Gossypoides kirkii</i> | <i>Malvaceae</i> |
| 10 | 11 | Eukaryota <i>Spondias tuberosa</i> | <i>Anacardiaceae</i> |

1

2 **Table S3:** BlastN screening statistics for Pacbio HiFi reads. Screening was performed against the NCBI
3 nucleotide database, Readouts are indicated in counts of *Malvaceae* species specific and non-*Malvaceae*
4 hits for a subset of 1000 HiFi reads.

5

| Rank | No. hits | Hits |
|--------------|-------------|------------------------------------|
| 1 | 889 (88.9%) | No hits |
| 2 | 64 (6.4%) | <i>Malvaceae</i> mitochondria |
| 3 | 47 (4.7%) | Non- <i>Malvaceae</i> mitochondria |
| Total | 111 (11.1%) | mitochondria |
| Rank | No. hits | Hits |
| 1 | 919 (91.9%) | No hits |
| 2 | 44 (4.4%) | Non- <i>Malvaceae</i> chloroplast |
| 3 | 37 (3.7%) | <i>Malvaceae</i> chloroplast |
| Total | 81 (8.1%) | chloroplast |

6

7 **Table S4:** Pacbio sequence library contamination statistics for 1000 HiFi. Organelle content was
8 determined using a BlastN screening against mitochondrial and chloroplast databases.

| Assembly statistic | Primary ctgs | Alternative ctgs | Hybrid scfds |
|--------------------|--------------|------------------|--------------|
| Ctgs/scfds | 1417 | 526 | 124 |
| Total length | 1,223.6 Mb | 17.2 Mb | 1,194.5 Mb |
| Median length | 32.4 kb | 24.5 kb | 16.320 Mb |
| Max length | 25.1 Gb | 652 kb | 29.444 Mb |
| Min length | 13.3 kb | 10.3 kb | 125 kb |
| N50 length | 10.6 Gb | 30.5 kb | 18.929 Mb |
| N50 index | 43 | 126 | 27 |
| N95 length | 483 kb | 15.5 kb | 7.206 Mb |
| N95 index | 168 | 465 | 64 |
| GC content | 34.36% | 47.6% | 33.76% |

1

2 **Table S5:** NGS assembly and hybrid scaffolding statistics. Sequences were assembled using the Hifiasm
3 assembler and scaffold with Bionano Genomics genome maps.

4










| Genome map statistic | Count |
|--|-----------|
| Genome map count | 216 |
| Label density (/100kb) | 15.94 |
| Total genome map length (Mbp) | 1,248.8 |
| Genome map N50 (Mbp) | 12.976 |
| Total molecules aligned to genome maps | 3,146,963 |
| Fraction of molecules aligned | 0.572 |
| Effective coverage | 374.509 |
| Average confidence | 21.8 |

5

6 **Table S6:** *De novo* genome map assembly statistics. Assembled molecules were mapped back to
7 genome maps to estimate the effective coverage and average confidence the *de novo* assembly.

8

9

| BUSCO class | | BUSCO distribution and topology | | | | | | |
|----------------|-------|---------------------------------|---|--------|--|--------|---|--------|
| Copy nr | Count | 1 ctg | Config. | 2 ctgs | Configuration | 3 ctgs | Configuration | 4 ctgs |
| Single | 284 | 284 |  | 0 | | 0 | | 0 |
| Duplicated | 1150 | 1 |  | 1149 |  | 0 | | 0 |
| Triuplicated | 843 | 0 | | 6 |  | 837 |  | 0 |
| Quadruplicated | 7 | 0 | | 1 |  | 6 |  | 0 |
| Quintuplicated | 3 | 0 | | 0 | | 3 |  | 0 |
| Sextuplicated | 0 | 0 | | 0 | | 0 | | 0 |
| Septuplicated | 1 | 0 | | 0 | | 1 |  | 0 |
| Total | 2288 | 285 | | 1156 | | 847 | | 0 |

1

2 **Table S7:** BUSCO distribution and topology. BUSCO genes are classified according to their copy number
3 in the genome. Distribution counts for ortholog gene copies have been indicated according to their
4 position either on one, two or three contigs. Configuration of gene copies is depicted by a horizontal line
5 representing a contig and superimposed small grey coloured boxes representing a gene copy.

| Species | n/x | ploidy | identified | complete | duplicate | fragmented | missing |
|--------------------|----------|----------------|------------|----------|-----------|------------|---------|
| A. thaliana | 2n=2x=10 | diploid | 98% | 98% | 17% | 0.5% | 2% |
| B. campestris | 2n=2x=18 | diploid | 99.4% | 99.1% | 14.7% | 0.3% | 0.6% |
| B. napus | 2n=4x=38 | allotetraploid | 97% | 97% | 90% | 1% | 3% |
| T. occidentale | 2n=2x=16 | diploid | 94% | 84% | 10% | 3% | 3% |
| T. pallescens | 2n=2x=16 | diploid | 94% | 83% | 11% | 2% | 4% |
| T. repens | 2n=4x=32 | allotetraploid | 94% | 24% | 57% | 3% | 6% |
| O. latifolia | 2n=2x=22 | diploid | 94.5% | 89% | 35% | 4.8% | 5.5% |
| B. amplexicaulis | 2n=6x=72 | hexaploid | 93.3% | 85% | 59% | 4.8% | 6.7% |
| S. lycopersicum | 2n=2x=24 | diploid | 96.8% | 96.4% | 1.0% | 0.3% | 2.1% |
| S. pennellii | 2n=2x=24 | diploid | 96.6% | 96.1% | 1.5% | 0.5% | 2.0% |
| S. lycopersicoides | 2n=2x=24 | diploid | 87.7% | 87.2% | 10.4% | 0,5% | 2.0% |

| | | | | | | | |
|-----------------|----------|----------------|-------|-------|-------|------|------|
| S. chacoense M6 | 2n=2x=24 | diploid | 97% | 96% | 4.3% | 1% | 3.0% |
| S. tuberosum RH | 2n=2x=24 | diploid | 98.6% | 97% | 74.1% | 1.6% | 1.4% |
| S. tuberosum | 2n=4x=48 | autotetraploid | 94.2% | 85.7% | 9.5% | 8.5% | 5.7% |

1

| Tissue | Library type | Polymerase reads | | | Subreads | | Insert | |
|--------|--------------|------------------|-----------|------------------|------------------|----------|------------------|----------|
| | | Size (Gbp) | Count | Mean length (bp) | Mean length (bp) | N50 (bp) | Mean length (bp) | N50 (bp) |
| Leaf | IsoSeq | 287.27 | 3,328,197 | 86,312 | 2,011 | 2,137 | 5,276 | 11,650 |
| Flower | IsoSeq | 323.48 | 6,676,796 | 48,448 | 1,723 | 1,774 | 4,536 | 8,917 |
| Pod | IsoSeq | 354.72 | 7,008,176 | 50,614 | 1,698 | 1,757 | 4,437 | 8,011 |

2

3 **Table S8:** Transcriptome sequencing statistics. Pacbio IsoSeq libraries were constructed for 3 different
4 tissues as indicated.

5

| Sample | Type | Size (Gbp) | Mapper | Variant Caller | Map rate | SNPs |
|--------------|--------|------------|----------|----------------|----------|------|
| Haploid okra | IsoSeq | 20 | Minimap2 | GATK | 99.81% | 1109 |
| Commercial | IsoSeq | 1.6 | Minimap2 | GATK | 95.14% | 8127 |

6

7 **Table S9:** Transcriptome mapping to the okra reference genome and SNP calls.

8

| Superscaffold | configuration | #units | Start position | End position |
|---------------|---------------|--------|----------------|--------------|
| 100090 | 18S-5.8S-28S | 24 | 13 | 289182 |
| 100111 | 18S-5.8S-28S | 54 | 8117 | 578985 |
| 100264 | 18S-5.8S-28S | 26 | 1 | 315695 |
| 67 | 18S-5.8S-28S | 32 | 2927 | 361934 |
| 74 | 18S-5.8S-28S | 5 | 20588136 | 20647257 |
| 8 | 18S-5.8S-28S | 8 | 1887219 | 2602099 |
| 28 | 18S-5.8S-28S | 11 | 8169258 | 8465035 |
| 34 | 5S | 8674 | 12382260 | 14843474 |

| | | | | |
|----|----|-----|----------|----------|
| 34 | 5S | 71 | 14843474 | 18230209 |
| 44 | 5S | 16 | 8621481 | 8625812 |
| 62 | 5S | 852 | 890588 | 1182219 |

1

2 **Table S10:** Ribosomal gene clusters in the okra genome. Ribosomal gene clusters are characterized by
3 unit configuration and number of tandemly arranged unit copies per superscaffold. Total lengths of
4 clustered units that can be derived from the start and end position.

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

1 **Materials and Methods**

2 *Chromosome analysis*

3 Plants of the Green Star F1 hybrid of okra (*Abelmoschus esculentus*) were grown in small for
4 collecting actively growing rootlets that appeared at the outside of the pot soil. The root tips were
5 pretreated with 8-hydroxyquinolin and then fixed in freshly prepared glacial acetic acid : ethanol 96%
6 (1:3) and one day later transferred to ethanol 70% for longer storage at 4 °C. Young flower buds were
7 collected from nurse fields in Kamphaeng Saen, Thailand, and directly fixed in acetic acid ethanol without
8 pretreatment. Microscopic preparations of root tip mitoses and pollen mother cells at meiotic stages were
9 prepared following pectolytic enzyme digestion of cell walls and acetic acid maceration and cell spreading
10 following the protocol of Kantama *et al.* (2017). Air-dried slides were stained in 300 nM 4',6-diamidino-2-
11 phenylindole (DAPI) in Vectashield (Vector Laboratories) and studied under a Zeiss fluorescence
12 microscope equipped with 1.4 N.A. objectives and appropriate epifluorescence filters for DAPI. The
13 captured images were optimized for best contrast and brightness in Adobe Photoshop, and slightly
14 sharpened with the Focus Magic (www.focusmagic.com) 2D deconvolution sharpening to remove
15 excessive blurring of the DAPI fluorescence (Kantama *et al.*, 2017).

16 *Bionano optical maps*

17 Sequence-specific labelling of approximately 700 ng genomic DNA from okra cv. Green
18 Star and subsequent backbone staining and DNA quantification for BioNano mapping was done using a
19 Direct Label Enzyme (DLE-1, CTTAAG) according to the manufacturer protocol 30206F BioNano Prep
20 Direct Label and Stain Protocol ([https://bionanogenomics.com/wp-content/uploads/2018/04/30206-
21 Bionano-Prep-Direct-Label-and-Stain-DLS-Protocol.pdf](https://bionanogenomics.com/wp-content/uploads/2018/04/30206-Bionano-Prep-Direct-Label-and-Stain-DLS-Protocol.pdf)). Chip loading and real-time analysis was carried
22 out on a BioNano Genomics Saphyr® analyser, using the green color channel on 3 flow cells, according
23 to the manufacturer system guide protocol 30143C ([https://bionanogenomics.com/wp-
24 content/uploads/2017/10/30143-Saphyr-System-User-Guide.pdf](https://bionanogenomics.com/wp-content/uploads/2017/10/30143-Saphyr-System-User-Guide.pdf)). Using the DLE-1 enzyme, 1.18 Tbp of
25 filtered DNA molecules with an average length of 215kb was produced, with a label density of
26 15.9/100kb and a molecule N50 of 207 kb. Subsequently, a *de novo* assembly was constructed using
27 Bionano Access™ (v.3.2.1) and the non-haplotype aware assembly program without extend and split but
28 with cutting of the complex multi-path regions (CMPR). Per the default settings, molecules < 150 Kbp
29 were removed before assembly. Next, a hybrid scaffolding of assembled sequence contigs was performed
30 with Bionano Genomics Solve (v.3.2.1) with a 375X-fold coverage for the DLE-1 molecules. Molecule
31 quality hybrid scaffold report were carried out using the BioNano Solve™ analysis pipeline
32 (<https://bionanogenomics.com/support-page/data-analysis-documentation/>).

1 *Pacbio HiFi, linked-read sequencing and de novo assembly*

2 We produced 3 Pacbio HiFi libraries using gDNA isolated from okra leaf tissue according
3 to the manufacturers protocol (<https://www.pacb.com>). HiFi reads of 15-20 kb were generated by
4 Circular Consensus Sequencing, using 6 SMRT cells, in total yielding 1,400 Gbp of sequence data.
5 Subsequent consensus calling was done using the pbccs v5.0.0 command line utility. HiFi reads were
6 defined as CCS reads having a minimum number of 3 passes and a mean read quality score of Q20.
7 Reads from different libraries were then combined into a single dataset for further analyses. Assembly of
8 HiFi Reads was done using hifiasm v0.12-r304 for coverages of ~20X, ~84X and ~95X (Cheng *et al.*,
9 2020). Primary contigs of the ~95X coverage assembly were scaffolded using BIONANO GENOMICS SOLVE
10 v3.6_09252020 and an optical *de novo* assembly. Solve scaffolded output was further scaffolded, in
11 contrast to the unscaffolded output, using Arcs v1.2.2 (<https://github.com/bcgsc/arcs>) and Links v1.8.7
12 (<https://github.com/bcgsc/LINKS>) based on the 10X genomics data that was mapped using LONGRANGER
13 v2.2.2. Scaffolds resulting from the final step were renamed to fit the naming scheme from the Bionano
14 scaffolding.

15 The 10X Genomics libraries were constructed with the Chromium™ Genome Reagent Kits v2
16 (10X Genomics®) according to the Chromium™ Genome v2 Protocol (CG00043) as described by the
17 manufacturer (<https://www.10xgenomics.com>). 10X Genomics libraries were sequenced on 2 separate
18 runs using the Illumina Novaseq6000 platform and S2 flow cells. Base calling and initial quality filtering
19 of raw sequencing data was done using BCL2FASTQ v2.20.0.422 using default settings. The Long Ranger
20 pipeline from 10X Genomics was used to process the 800 Gbp sequencing output and align the reads to
21 the tomato reference genome. After detecting the conflict region with the BIONANO GENOMICS ACCESS SUITE
22 (v.1.3.0), we manually inspected the conflict regions using 10X linked-reads mapped to the
23 superscaffolds. Mapping and visualization of scaffolds was done with LONGRANGER WGA v.2.2.2 and LOUPE
24 v.2.1.1 respectively.

25 *Assembly QC*

26 Statistics on NX lengths, GC-percentage, mean-, median-, maximum- and minimum lengths of
27 contigs or scaffolds for each assembly step was collected from software output and when not available
28 was generated with custom python scripts. HiFi reads were mapped back to the assembly using MINIMAP2
29 v2.17-r941 to assess purging correctness with purge_haplotigs (Li, 2018; Roach *et al.*, 2018). MINIMAP2
30 alignment was also used for BLOOTOOLS v1.1.1 analysis to check the taxonomic origin, coverage and GC-
31 percentage of unscaffolded output (Laetsch *et al.*, 2017). The completeness of the assembly was
32 benchmarked using BUSCO with the eudicots_odb10 (eukaryota, 2020-09-10) lineage set to scan for
33 single copy orthologs (Kriventseva *et al.*, 2019; Simão *et al.*, 2015). AUGUSTUS v3.2.2 was subsequently

1 used for gene prediction. BUSCO output was used for topology analysis of duplicated genes (Stanke *et*
2 *al.*, 2006). To assess repeat content and synteny within the scaffolds, the assembly was self-aligned
3 using a combination of MINIMAP2 and DGENIES, and NUCMER v4.0.0beta2 together with MUMMERPLOT v3.5
4 (Kurtz *et al.*, 2004).

5 *Iso-Seq sequencing and data analysis*

6 Total RNA was isolated from leaf (10 µg), flower buds (21 µg) and young fruits (31µg) from okra
7 cv. "Green Star". RNA quality was checked on a Bioanalyzer platform (<https://www.agilent.com>) by
8 comparing to standard samples of 25S and 18S ribosomal RNA. Transcript samples were subsequently
9 used for construction of 3 sequence libraries and sequenced with PacBio SMRT technology
10 (<https://www.pacb.com/smrt-science/smrt-sequencing>) using 4 SMRT cells. Consensus reads were
11 produced with the CCS v5.0.0 command-line utility of PacBio. HiFi reads were classified as such using the
12 same specifications as the genomic reads. Primers sequences from reads were removed and demultiplexed
13 with LIMA v2.0.0. Poly-A tails were trimmed and concatemer were removed with ISOSEQ3 v3.4.0 refine to
14 generate full length non-concatemer reads and subsequently clustered with ISOSEQ3 CLUSTER. Since
15 distributions of mean read quality showed over 90% of data to have a mean quality score in range of 90-
16 93, no final polishing was applied. High quality full-length transcripts obtained from the SMRT analysis
17 pipeline were then mapped to the hybrid assembly using GMAP (Wu and Watanabe, 2005).

18 *Read analysis*

19 Quality control of reads was performed using SMRTLINK v9 and FASTQC v0.11.9. A random sample
20 of 1000 reads per SMRT cell was taken using SEQTK v1.3-r106 seq with a random seed from the BASH
21 v4.2.46 internal pseudorandom generator. Samples were screened for chloroplast content, plastid content
22 and taxonomic contamination applying BLAST v2.10.1+ with parameter settings *-evaluate 0.001* and *-*
23 *max_target_seqs 1* (Altschul *et al.*, 1990; Altschul *et al.*, 1997; Camacho *et al.*, 2009). The databases
24 used for each screening were NCBI nt, plastid and mitochondrion publicly available FTP downloads dated
25 2020-11-15 (Agarwala *et al.*, 2016).

26 K-mers were counted for both 10X linked reads and HiFi reads, using KMC 3.1.1 (Kokot *et al.*,
27 2017) with parameter settings *-m64 -ci1 -cs10000* for *k = 16, 21, 28, 37, 48* and *61* to determine the k-
28 mer size for best-model fit. For 10X linked reads 23 bp of R1, containing 16 bp barcode plus the 7 bp long
29 spacer sequence, were trimmed off before counting. The kmc histogram of counts was subsequently used
30 for estimation of genome parameters and visualization of k-mer spectrum using GENOMESCOPE2 (Ranallo-
31 Benavidez *et al.*, 2019). The polyploid nature of okra was further examined applying a locally developed,
32 publicly available fork of SMUDGE PLOT labeled v0.2.3dev_rn that is true to the original algorithm, allowing

1 for parallelization. In the original algorithm k-mer pairs with a hamming distance of 1 are found by a
2 recursive method that iterates over all positions within the k-mer. The redesigned algorithm parallelizes
3 the search by each thread looking at a given position within the k-mer. To reduce the number of false
4 negatives, results are then filtered using a bloom filter with an *error rate* set at 0.0001.

5 *Annotation*

6 Repeats annotation was carried out with REPEATMODELER, REPEATCLASSIFIER and REPEATMASKER tools
7 and the combined REPBASE (2014) and DFAM (2020) databases for classification of repeats (Bao *et al.* 2015;
8 Hubley *et al.*, 2016; Smith *et al.*, 2013). Full length non-concatemer IsoSeq reads were mapped against
9 the genome assembly using MINIMAP2 with parameter settings *-ax splice -uf --secondary=no -C5*. A repeat
10 masked genome and transcriptome read mapping was then input to the BRAKER2 v2.1.5 pipeline to
11 generate a structural annotation of genes using *ab initio* prediction (Borodovsky and Lomsadze, 2011; Hoff
12 *et al.*, 2019). Using the general feature format (GFF) file output, open reading frame translations of the
13 predicted genes were made with the default eukaryotic translation table. Produced polypeptide sequences
14 were then annotated with INTERPROSCAN v5.39-77 10 that collected annotations from the default set of
15 databases including PANTHER, GEN3D, CDD, etc (Hunter *et al.*, 2009; Jones *et al.*, 2014).

16 *Variant calling*

17 To investigate the diversity among okra accessions, multiple publicly available transcriptome
18 datasets were mapped to the assembled okra reference as presented in this study. Public datasets include
19 SRR620228 (<https://www.ncbi.nlm.nih.gov/sra/?term=SRR620228>), PRJNA393599
20 (<https://www.ncbi.nlm.nih.gov/sra/?term=PRJNA393599>) and PRJNA430490
21 (<https://www.ncbi.nlm.nih.gov/sra/?term=PRJNA430490>). In case multiple tissues or runs were available,
22 data were merged. Mapping was done using STAR v2.7.8a (Dobin *et al.*, 2013; Leinonen *et al.*, 2011).
23 Duplicates were marked applying GATK v4.2.0 MarkDuplicatesSpark mode (McKenna *et al.*, 2010). Then
24 variants were called using GATK HAPLOTYPECALLER with default filters. Obtained variants were then filtered
25 for SNPs and subsequently quality filtered by GATK SELECTVARIANTS with parameter filter settings *QD > 2.0*
26 *&& MQ > 40.0 && FS < 60.0 && SOR < 3.0 && QUAL > 30.0 && MQRankSum > -12.5 && ReadPosRankSum*
27 *> -8.0*. IsoSeq reads were mapped as described above.

28 *Subgenome separation*

29 To separate the subgenomes in the okra genome assembly, a clustering approach on the repetitive
30 part of the okra genome was applied to generate clusters of scaffolds with similar repeat patterns. For each
31 scaffold, forward and reverse complement non-canonical k-mers were counted setting a length *k = 13*.
32 The counts were normalized by scaffold length and then filtered to have a minimal count of 100. Sets of

1 k-mer counts for both forward and reverse strand of each scaffold were then merged. All counts were
2 subsequently increased by one and log₁₀ transformed to generate a scale that was appropriate for
3 visualization. Euclidean distances between the scaffolds were stored and scaffolds were subsequently
4 clustered using the Ward's (minimum variance) method. To visualize the clustering along the k-mer
5 patterns, we generated a cluster map from a heatmap for the k-mer counts combined with a hierarchical
6 clustering of the scaffolds. In addition, the clustering of kmers generated from an in-house generated
7 artificial potato tomato hybrid assembly, and the allopolyploid cotton genome was used as a test set.
8 Clusters were then compared based on BUSCO score completeness and scanned for homoeology in a
9 network visualization using an adjacency matrix. Scores in the adjacency matrix were taken from the edge
10 counts between scaffolds that have identical BUSCO genes in common, implying that scaffold x and scaffold
11 y have an edge count z when sharing a number of z BUSCO genes. *Statistical analysis of metabolic pathway*
12 *assignment for okra orthologs*

13 The mapping probability for okra orthologs to KEGG reference pathways was based on a
14 hypergeometric test (one-sided Fisher's exact test) to measure the statistical significance for pathway
15 assignment of a putative okra ortholog set. Pathway *p*-values were calculated according to equation 1
16 (Eq. 1), where K equals the unique enzymes known for a pathway *p*, *k* for the number of searched
17 enzymes uniquely mapping on pathway *p*, N as the number of unique enzymes of all reference species
18 known for all pathways, and *n* the number of searched enzymes uniquely mapping on all pathways.

19 (Eq. 1)
$$P(X = k) = f(k, N, K, n) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}$$

20

21

22

23

24

25

26

27

28

1 **Literature**

- 2 **Aflitos, S., Schijlen, E., de Jong, H. et al.** (2014) Exploring genetic variation in the tomato (*Solanum*
3 section *Lycopersicon*) clade by whole-genome sequencing. *Plant J.* **80**, 136–148.
- 4 **Ali, S., Khan, M.A., Rasheed, H.S. and Iftikhar, Y.** (2005) Management of Yellow Vein Mosaic disease
5 of Okra through pesticide/bio-pesticide and suitable cultivars. *Int. J. Agric. Biol.* **7**, 145-147.
- 6 **Agarwala, R., Barrett, T., Beck, J., Benson, D. A., Bollin, C., Bolton, E. et al.** (2016). Database
7 resources of the National Center for Biotechnology Information. *NAR*, **44** (Database issue), D7.
- 8 **Altschul, S. F., Gish, W., Miller, W., Myers, E. W. and Lipman, D. J.** (1990). Basic local alignment
9 search tool. *J. Mol. Biol.*, **215**, 403–410.
- 10 **Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W., & Lipman, D. J.**
11 (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *NAR* **25**,
12 3389–3402.
- 13 **Baird, D.M.** (2017) Telomeres and genomic evolution. *Phil. Trans. R. Soc. B* **373**, 20160473.
- 14 **Bao, W., Kojima, K. K. and Kohany, O.** (2015). Repbase Update, a database of repetitive elements in
15 eukaryotic genomes. *Mobile DNA* **6**, 1–6.
- 16 **Benchasri, S.** (2012) Okra (*Abelmoschus esculentus* (L.) Moench) as a valuable vegetable of the world.
17 *Ratar. Povrt.* **49**, 105-112.
- 18 **Bird, K.A., VanBuren, R., Puzey, J.R. and Edger, P.P.** (2018) The causes and consequences of
19 subgenome dominance in hybrids and recent polyploids. *New Phyt.* **220**, 87-93.
- 20 **Borodovsky, M. and Lomsadze, A.** (2011). Eukaryotic gene prediction using GeneMark.hmm-E and
21 GeneMark-ES. *Curr. Protoc. Bioinformatics* **4**, 610.
- 22 **Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. and Madden, T. L.**
23 (2009). BLAST+: Architecture and applications. *BMC Bioinformatics* **10**, 1–9.
- 24 **Cheng, H., Concepcion, G.T., Feng, X., Zhang, H. and Li, H.** (2021) Haplotype-resolved *de novo*
25 assembly using phased assembly graphs with hifiasm. *Nature Methods* **18**, 170-175.
- 26 **Dankhar, S.K., and Koundinya, A.V.V.** (2020) Accelerated breeding in Okra. In *Accelerated plant*
27 *breeding, volume 2. Vegetable crops* (Gosal, S.S., and Shabir Hussain Wani, S.H., eds). Springer Nature
28 Switzerland AG, Switzerland, pp 337-354.

- 1 **Dunwell, J.M.** (2010) Haploids in flowering plants: origins and exploitation. *Plant Biotech. J.* **8**, 377-
2 424.
- 3 **Goffová, I. and Fajkus, J.** (2021) The rDNA loci-Intersections of replication, transcription, and repair
4 pathways. *MDPI* **22**, 1302.
- 5 **Griffiths, A.G., Moraga, R., Tausen, M., Gupta, V., Bilton, T.P. et al.** (2019). Breaking free: The
6 genomics of allopolyploidy-facilitated niche expansion in white clover. *Plant Cell* **31**, 1466-1487.
- 7 **Guo, Z-H., Ma, P-F., Yang, G-Q., Hu, J-Y., Liu, Y-L. et al.** (2019) Genome sequence provides insights
8 into the reticulate origin and unique traits of woody bamboos. *Mol. Plant* **12**, 1353-1365.
- 9 **Hoff K.J., Lomsadze A., Borodovsky M. and Stanke M.** (2019) Whole-Genome Annotation with
10 BRAKER. In: Kollmar M. (eds) Gene Prediction. Methods in Molecular Biology, vol 1962. Humana, New
11 York, NY.
- 12 **Hubley, R., Finn, R. D., Clements, J., Eddy, S. R., Jones, T. A., Bao, W. et al.** (2016). The
13 Dfam database of repetitive DNA families. *NAR* **44**, D81.
- 14 **Hu, Y., Chen, J., Fang, L. et al.** (2019) *Gossypium barbadense* and *Gossypium hirsutum* genomes
15 provide insights into the origin and evolution of allotetraploid cotton. *Nat Genet* **51**, 739–748.
- 16 **Hunter, S., Apweiler, R., Attwood, T. K., Bairoch, A., Bateman, A., Binns, D. et al.** (2009).
17 InterPro: The integrative protein signature database. *NAR* **37**, D211-D215.
- 18 **Jones, P., Binns, D., Chang, H. Y., Fraser, M., Li, W., McAnulla, C. et al.** (2014). InterProScan 5:
19 Genome-scale protein function classification. *Bioinformatics* **30**, 1236-1240.
- 20 **Joshi A. B. and Hardas M. W.** (1956) Allopolyploid nature of okra, *Abelmoschus esculentus* (L.) Moench.
21 *Nature* **178**, 1190.
- 22 **Kantama, L., Wijnker, E. and de Jong, H.** (2017) Optimization of Cell Spreading and Image Quality
23 for the Study of Chromosomes in Plant Tissues. In: Plant Germline Development (Schmidt, A. ed).
24 Methods in Molecular Biology, vol 1669. Humana Press, New York, NY., pp 141-158.
- 25 **Kokot, M., Dlugosz, M. and Deorowicz, S.** (2017). KMC 3: counting and manipulating k-mer statistics.
26 *Bioinformatics* **33**, 2759–2761.

- 1 **Kriventseva, E. V., Kuznetsov, D., Tegenfeldt, F., Manni, M., Dias, R., Simão, F. A. and Zdobnov,**
2 **E. M.** (2019). OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral
3 genomes for evolutionary and functional annotations of orthologs. *NAR*, 47(D1), D807–D811.
- 4 **Kumar, S., Jones, M., Koutsovoulos, G., Clarke, M. and Blaxter, M.** (2013) Blobology: exploring
5 raw genome data for contaminants, symbionts and parasites using taxon-annotated GC-coverage plots.
6 *Front. Genet.* **4**, 237.
- 7 **Kurtz, S., Phillippy, A., Delcher, A. L., Smoot, M., Shumway, M., Antonescu, C. and Salzberg, S.**
8 **L.** (2004). Versatile and open software for comparing large genomes. *Gen. Biol.* **5**, 1–9.
- 9 **Kyriakidou, M., Anglin, N.L., Ellis, D., Tai, H.H. and Strömviik, M.V.** (2020). Genome assembly of
10 six polyploid potato genomes. *Sci. Data*, **7**, 88.
- 11 **Laetsch, D. R., Blaxter, M. L., Eren, A. M., and Leggett, R. M.** (2017). BlobTools: Interrogation of
12 genome assemblies. *F1000Research* **6**, 1287.
- 13 **Langley, C.H., Crepeau, M., Cardeno, M., Corbett-Detig, R. and Stevens, K.** (2011) Circumventing
14 heterozygosity: sequencing the amplified genome of a single haploid *Drosophila melanogaster* embryo.
15 *Genetics* **188**, 239-246.
- 16 **Lata, S., Yadav, R.K. and B.S. Tomar, B.S.** (2021). Genomic tools to accelerate improvement in okra
17 (*Abelmoschus esculentus*), Landraces - Traditional Variety and Natural Breed, Amr Elkelish, IntechOpen,
18 DOI: 10.5772/intechopen.97005.
- 19 **Li, F., Fan, G., Lu, C., Xiao, G., Zou, C., Kohel, R. J. et al.** (2015). Genome sequence of cultivated
20 Upland cotton (*Gossypium hirsutum* TM-1) provides insights into genome evolution. *Nat. Biotech.* **33**,
21 524–530.
- 22 **Li, J., Ye, G-y., Liu, H-l. and Wang, Z-h.** (2020) Complete chloroplast genomes of three
23 important species, *Abelmoschus moschatus*, *A. manihot* and *A. sagittifolius*: Genome structures,
24 mutational hotspots, comparative and phylogenetic analysis in *Malvaceae*. *PLoS One* **15**, e0242591.
- 25 **Liu, I.M., Liou, S.S., Lan, T.W., Hsu, F.L., Cheng, J.T.** (2005) Myricetin as the active principle of
26 *Abelmoschus moschatus* to lower plasma glucose in streptozotocin-induced diabetic rats. *Planta Med.* **71**,
27 617–21.

- 1 **Merita, K., Kattakunnel, J.J., Yadav, S.R., Bhat, K.V. and Rao, S.R.** (2012) Chromosome counts in
2 wild and cultivated species of *Abelmoschus Medikus*. from the Indian sub-continent. *J. Hort Sci. Biot.* **87**,
3 593-599.
- 4 **Moriya, Y., Itoh, M., Okuda, S., Yoshizawa, A.C. and Kanehisa, M.** (2007) KAAS: an automatic
5 genome annotation and pathway reconstruction server. *Nucleic Acids Research* **35**, 182–185.
- 6 **Muimba-Kankolonga, A.** (2018) Vegetable production. In *Food crop production by smallholder farmers*
7 *in Southern Africa* (Demetre, C., ed). Academic Press, London, pp. 205-273.
- 8 **Naumova, T.N.** (2008) Apomixis and amphimixis in flowering plants. *Cyt. Genet.* **42**, 53-65.
- 9 **Novák, P., Guignard, M. S., Neumann, P., Kelly, L. J., Mlinarec, J., Koblížková, A. et al.** (2020)
10 Repeat-sequence turnover shifts fundamentally in species with large genomes. *Nat. Plants* **6**, 1325–
11 1329.
- 12 **Peter, E.L., Nagendrappa, P.B., Ajayi, C.O., Sesaazi, C.D.** (2021) Total polyphenols and
13 antihyperglycemic activity of aqueous fruits extract of *Abelmoschus esculentus*: Modeling and
14 optimization of extraction conditions. *PLoS ONE* **16**, e0250405.
- 15 **Peters, S.A., Datema E., Szinay, D. et al.** (2009) *Solanum lycopersicum* cv. Heinz 1706 chromosome
16 6: distribution and abundance of genes and retrotransposable elements. *Plant J.*, **58**, 867-869.
- 17 **Portemer, V., Renne, C., Guillebaux, A. and Mercier, R.** (2015) Large genetic screens for
18 gynogenesis and androgenesis haploid inducers in *Arabidopsis thaliana* failed to identify mutants. *Front.*
19 *Plant Sci.* **6**, 581–6.
- 20 **Ranallo-Benavidez, T.R., Jaron, K.S. and Schatz, M.C.** (2020). Genomescope 2.0 and smudgeplot for
21 reference free profiling of polyploid genomes. *Nature Comm.* **11**, 1432.
- 22 **Roach, M. J., Schmidt, S. A. and Borneman, A. R. (2018).** Purge Haplotigs: allelic contig
23 reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* **19**, 460.
- 24 **Salameh, N.** (2014) Flow cytometric analysis of nuclear DNA of Okra Landraces (*Abelmoschus*
25 *esculentus* L.). *Am. J. Agric. Biol. Sci.* **9**, 245-250.
- 26 **Siemonsma, J.S.** (1982) West African Okra - Morphological and cytogenetical indications for the
27 existence of a natural amphidiploid of *Abelmoschus esculentus* (L.) Moench and *A. Manihot* (L.) Medikus.
28 *Euphytica* **31**, 241-252.

- 1 **Simaõ , F.A., Waterhouse, R.M., Ioannidis, P., Evgenia V. Kriventseva, E.V. and Zdobnov, E.M.**
2 (2015) BUSCO: assessing genome assembly and annotation completeness with single copy-orthologs.
3 *Bioinformatics* **31**, 3210-3212.
- 4 **Smith, A., Hubley, R. and Green, P.** (2013). RepeatMasker Open-4.0. *RepeatMasker Open-4.0*.
- 5 Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S. and Morgenstern, B. (2006). AUGUSTUS: *ab*
6 *initio* prediction of alternative transcripts. *NAR* **34**, W435–W439.
- 7 **Takakura, K-I, and Nishio, T.** (2012) Safer DNA extraction from plant tissues using sucrose buffer and
8 glass fiber filter. *J. Plant Res.* **125**, 805-807.
- 9 **The Tomato Genome Consortium** (2012) The tomato genome sequence provides insights into fleshy
10 fruit tomato. *Nature* **485**, 635–641.
- 11 **Venkataravanappa, V., Lakshminarayana Reddy, C.N. and Krishna Reddy, M.** (2013)
12 Begomovirus characterization, and development of phenotypic and DNA-based diagnostics for screening
13 of okra genotype resistance against *Bhendi yellow vein mosaic virus*. *3 Biotech* **3**, 461-470.
- 14 **Wu, D-T., Nie, X-R., Li, H-Y., et al.** (2020) Phenolic compounds, antioxidant activities, and inhibitory
15 effects on digestive enzymes of different cultivars of okra (*Abelmoschus esculentus*). *MDPI* **25**, 1276.
- 16 **Wu, T.D. and Watanabe, C.K.** (2005) GMAP: a genomic mapping and alignment program for mRNA
17 and EST sequences. *Bioinformatics*, **21**, 1859–1875.
- 18