1  # High contiguity *de novo* genome sequence assembly

2  # of Trifoliate yam (*Dioscorea dumetorum*) using long

3  # read sequencing

4

5  **Christian Siadjeu[1,2,+], Boas Pucker[2,3,+], Prisca Viehöver[2], Dirk C. Albach[1] and Bernd Weisshaar[2]**[*]

6  [1]  Institute for Biology and Environmental Sciences, Biodiversity and Evolution of Plants, Carl-von-Ossietzky
7      University Oldenburg, Carl-von-Ossietzky Str. 9-11, 26111 Oldenburg, Germany; christian.siadjeu@uol.de
8      (CS); dirk.albach@uol.de (DCA)
9  [2]  Genetics and Genomics of Plants, Faculty of Biology, Center for Biotechnology (CeBiTec), Bielefeld
10     University, Sequenz 1, 33615 Bielefeld, NRW, Germany; bpucker@cebitec.uni-bielefeld.de (BP),
11     viehoeve@cebitec.uni-bielefeld.de (PV), bernd.weisshaar@uni-bielefeld.de (BW)
12  [3]  Molecular Genetics and Physiology of Plants, Faculty of Biology and Biotechnology, Ruhr-University
13     Bochum, Universitätsstraße 150, 44801 Bochum, Germany; bpucker@cebitec.uni-bielefeld.de (BP)
14  *  Correspondence: bernd.weisshaar@uni-bielefeld.de; Tel.: +49-521-106-8720
15  +  Shared first authorship. CS and BP contributed equally to this work and are co-first authors.

16  **Abstract:** Trifoliate yam (*Dioscorea dumetorum*) is one example of an orphan crop, not traded
17  internationally. Post-harvest hardening of the tubers of this species starts within 24 hours after
18  harvesting and renders the tubers inedible. Genomic resources are required for *D. dumetorum* to
19  improve breeding for non-hardening varieties as well as for other traits. We sequenced the *D.*
20  *dumetorum* genome and generated the corresponding annotation. The two haplophases of this
21  highly heterozygous genome were separated to a large extent. The assembly represents 485 Mbp of
22  the genome with an N50 of over 3.2 Mbp. A total of 35,269 protein-encoding gene models as well as
23  9,941 non-coding RNA genes were predicted and functional annotations were assigned.

24

25  **Keywords:** yam; *D. dumetorum*; nanopore sequencing; genome assembly; comparative genomics;
26  read depth

27

28  ## 1. Introduction

29     The yam species *Dioscorea dumetorum* (trifoliate yam) belongs to the genus *Dioscorea* comprising
30  about 600 described species. The genus is widely distributed throughout the tropics [1] and includes
31  important root crops that offer staple food for over 300 million people. Eight *Dioscorea* species are
32  commonly consumed in West and Central Africa, of which *D. dumetorum* has the highest nutrient
33  value [2]. Tubers of *D. dumetorum* are protein-rich (9.6%) with a fairly balanced essential amino acids
34  composition [3]. The provitamin A and carotenoid contents of the tubers of deep yellow genotypes
35  are equivalent to those of yellow corn maize lines selected for increased concentrations of
36  provitamin A [4]. The deep yellow yam tubers are used in antidiabetic treatments in Nigeria [5],
37  probably due to the presence of dioscoretine, which is a bioactive compound with hypoglycaemic
38  properties [6]. Yet, *D. dumetorum* constitutes an underutilized and neglected crop species despite its
39  great potential for nutritional, agricultural, and pharmaceutical purposes.
40     Unlike other yam species, the agricultural value of *D. dumetorum* is limited by post-harvest
41  hardening, which starts within 24 h after harvest and renders tubers inedible. Previous research
42  showed that among 32 *D. dumetorum* cultivars tested, one cultivar was not affected by the hardening
43  phenomenon [7]. This discovery provides a starting point for a breeding program of *D. dumetorum*

44  against the post-harvest hardening phenomenon. *Dioscorea* cultivars are obligate outcrossing plants
45  that display highly heterozygous genomes. Thus, methods of genetic analysis routinely used in
46  inbreeding species such as linkage analysis using the segregation progeny of an F2 generation and
47  recombinant inbred lines are inapplicable to yam [8]. Furthermore, the development of
48  marker-assisted selection requires the establishment of marker assays and dense genetic linkage
49  maps. Thus, access to a complete and well-annotated genome sequence is one essential step towards
50  the implementation of comprehensive genetic, genomic and population genomics approaches for *D.*
51  *dumetorum* breeding. So far, a genome sequence assembly for *D. rotundata* (Guinea yam) [8] and a
52  reference genetic map for *D. alata* (Greater yam) [9] have been released. However, these two species
53  belong to the same section of *Dioscorea* (*D.* sect. *Enantiophyllum*) but are distant from *D. dumetorum*
54  (*D.* sect. *Lasiophyton*) in phylogenetic analyses [10,11]. They also differ in chromosome number
55  [8,12,13] making it unlikely that genetic maps can be directly transferred to *D. dumetorum*. Here, we
56  report long read sequencing and *de novo* genome sequence assembly of the *D. dumetorum* Ibo sweet 3
57  cultivar that does not display post-harvest hardening.
58

59  **2. Materials and Methods**

60  *2.1. Sampling and Sequencing*

61  The *D. dumetorum* accession Ibo sweet 3 that does not display post-harvest hardening had been
62  collected in the South-West region of Cameroon in 2013 [7]. Tubers of this accession were transferred
63  to Oldenburg (Germany) and the corresponding plants were cultivated in a greenhouse at 25°C. The
64  haploid genome size of the Ibo sweet 3 genotype had been estimated to be 322 Mbp through flow
65  cytometry [14].
66  High molecular weigth DNA was extracted from 1g of leaf tissue using a CTAB-based method
67  modified from [15]. After grinding the sample in liquid nitrogen, the powder was suspended in 5
68  mL CTAB1 (100 mM Tris-HCl pH 8.0, 20 mM EDTA, 1.4 M NaCl, 2% CTAB, 0.25% PVP) buffer
69  supplemented with 300 µL ß-mercaptoethanol. The suspension was incubated at 75°C for 30
70  minutes and inverted every five minutes. Next, 5 mL dichloromethane were added and the solutions
71  were mixed by inverting. The sample was centrifuged at 11,200 g at 20°C for 30 minutes. The clear
72  supernatant was mixed with 10 mL CTAB2 (50 mM Tris-HCl pH 8.0, 10 mM EDTA, 1% CTAB,
73  0.125% PVP) in a new reaction tube by inverting. Next, a centrifugation was performed at 11,200 g at
74  20°C for 30 minutes. After discarding the supernatant, 1 mL NaCl (1 M) was added to re-suspend
75  the sediment by gently flicking the tube. By adding an equivalent amount of 1mL isopropanol and
76  careful mixing, the DNA was precipitated again and the sample was centrifuged as described above.
77  After washing the sediment with 1 mL of 70% ethanol, 200 µL TE buffer (10 mM Tris pH 8.0, 0.1 mM
78  EDTA) containing 2 mg DNAse-free RNaseA were added. Re-suspension and RNA degradation
79  were achieved by incubation over night at room temperature. DNA quality and quantity were
80  assessed via NanoDrop2000 measurement, agarose gel electrophoresis, and Qubit measurement.
81  The short read eliminator (SRE) kit (Circulomics) was used to enrich long DNA fragments following
82  the suppliers' instructions. Results were validated via Qubit measurement.
83  Library preparation was performed with 1 µg of high molecular weight DNA following the
84  SQK-LSK109 protocol (Oxford Nanopore Technologies, ONT). Sequencing was performed on four
85  R9.4.1 flow cells on a GridION. Flow cells were treated with nuclease flush (20 µL DNaseI (NEB) and
86  380 µL nuclease flush buffer) once the number of active pores dropped below 200, to allow
87  successive sequencing of multiple libraries on an individual flow cell. Live base calling was
88  performed on the GridION by Guppy v3.0 (ONT).
89  A total of 200 ng high molecular weight gDNA was fragmented by sonication using a Bioruptor
90  (Diagenode) and subsequently used for Illumina library preparation. End-repaired fragments were
91  size selected by AmpureXp Beads (Beckmann-Coulter) to an average size of 650 bp. After A-tailing
92  and adaptor ligation fragments that carry adaptors on both ends were enriched by 8 cycles of PCR
93  (Illumina TruSeq Nano DNA Sample Kit). The final library was quantified using PicoGreen

94　(Quant-iT) on a FLUOstar plate reader (BMG labtech) and quality checked by HS-Chips on a 2100
95　Bioanalyzer (Agilent Technologies). The PE library was sequenced in 2 x 250 nt mode on an Illumina
96　HiSeq-1500.
97

98　*2.2. Genome assembly and polishing*

99　Genome size prediction was performed with GenomeScope [16], findGSE [17], and gce [18]
100　based on k-mer histograms generated by JellyFish v2 [19] as previously described [20] for different
101　k-mer size values. In addition, MGSE [20] was run on an Illumina read mapping with single copy
102　BUSCOs as reference regions for the haploid coverage calculation. Smudgeplot [21] was run on the
103　same k-mer histograms (also for different k-mer size values) as the genome size estimations to
104　estimate the ploidy.
105　Canu v1.8 [22] was deployed for the genome assembly. Raw ONT reads were provided as input
106　to Canu for correction and trimming. Subsequently, Canu assembled the genome sequence from the
107　resulting polished reads. The following optimized parameters were used "'genomeSize = 350m',
108　'corOutCoverage = 200' 'correctedErrorRate = 0.12' batOptions = -dg 3 -db 3 -dr 1 -ca 500 -cp 50'
109　'minReadLength = 10000' 'minOverlapLength = 5000' 'corMhapFilterThreshold = 0.0000000002'
110　'ovlMerThreshold = 500' 'corMhapOptions = --threshold 0.85 –num-hashes 512 –num-min-matches 3
111　–ordered-sketch-size 1000 –ordered-kmer-size 14 –min-olap-length 5000 –repeat-idf-scale 50'". The
112　parameters we selected were optimized for the assembly of a heterozygous genome sequence and
113　our data set. The value for the genome size, estimated to be 322 Mbp, was increased to 350 Mbp to
114　increase the number of reads utilized for the assembly process. A total of 66.7 Gbp of ONT reads
115　with an N50 of 23 kbp was used for assembly, correction and trimming.
116　ONT reads were mapped back to the assembled sequence with minimap2 v2.17 [23], using the
117　settings recommended for ONT reads. Next, the contigs were polished with racon v.1.4.7 [24] with
118　-m 8 -x -6 -g -8 as recommended prior to the polishing step with medaka. Two runs of medaka
119　v.0.10.0 (https://github.com/nanoporetech/medaka) polishing were performed with default
120　parameters (-m r941_min_high) using ONT reads. Illumina short reads were aligned to the medaka
121　consensus sequence using BWA-MEM v. 0.7.17 [25]. This alignment was subjected to Pilon v1.23 [26]
122　for final polishing in three iterative rounds with default parameters for the correction of all variant
123　types and –mindepth 4.
124　Downstream processing was based on a previously described workflow [27] and performed by
125　customized Python scripts for purging of contigs shorter than 100 kbp and calculation of assembly
126　statistics (https://github.com/bpucker/yam). In general, sequences were kept if matching a white list
127　(*D. rotundata*) and discarded if matching a black list (bacterial/fungal genome sequences). Sequences
128　with perfect matches against the genome sequences of plants that were sequenced in the lab in
129　parallel (*A. thaliana*, *Beta vulgaris*, and *Vitis vinifera*) were discarded as well. Contigs with less than
130　3-fold average coverage in an Illumina short read mapping were compared against nt via BLASTn
131　with an e-value cut-off at $10^{-10}$ to identify and remove additional bacterial and fungal sequences.
132　For the ordering ("scaffolding" according to linkage groups) the *D. dumetorum* assembly we
133　employed *D. rotundata* pseudochromosomes. *D. rotundata* pseudochromosome sequences longer
134　than 100 kbp were split into chunks of 1000 bp and subject to a BLASTn search against the *D.
135　dumetorum* assembly with a word size of 12. Hits were considered if the similarity was at least 70%
136　and if at least 70% of the query length were covered by the alignment. To avoid ambiguous hits
137　against close paralogs or between repeat units, BLAST hits were exclude if the second hit exceeds
138　90% of the score of the top hit. The known order of all chunks on the *D. rotundata* sequence was
139　considered as a "pseudo genetic map" to arrange the *D. dumetorum* contigs via ALLMAPS v0.9.14
140　[28].

141　*2.3. Genome sequence annotation*

142　Hints for gene prediction were generated by aligning *D. rotundata* transcript sequences (TDr96
143　v1.0) [8] as previously described [29]. BUSCO v3 [30] was applied to generate a species-specific set of

144  AUGUSTUS gene prediction parameter files. For comparison of annotation results, the *D. rotundata*
145  genome assembly GCA_002260605.1 [8] was retrieved from NCBI. Gene prediction hints of *D.*
146  *dumetorum* and dedicated parameters were subjected to AUGUSTUS v.3.3 [31] for gene prediction
147  with previously described settings [29]. Various approaches involving AUGUSTUS parameter files
148  for rice and maize genome sequences provided by AUGUSTUS as well as running the gene
149  prediction on a sequence with repeats masked by RepeatMasker v4.0.8 [32] with default parameters
150  were evaluated. BUSCO was applied repeatedly to assess the completeness of the gene predictions.
151  The best results for *D. dumetorum* genome sequence annotation were obtained by using an
152  unmasked assembly sequence and by applying yam specific AUGUSTUS gene prediction parameter
153  files generated via BUSCO as previously described [30,33]. Predicted genes were filtered based on
154  sequence similarity to entries in several databases (UniProt/SwissProt, Araport11, *Brachypodium*
155  *distachyon* v3.0, *Elaeis guineensis* v5.1, GCF_000005425.2, GCF_000413155.1, *Musa acuminata* Pahang
156  v2). Predicted peptide sequences were compared to these databases via BLASTp [34] using an
157  e-value cut-off of $10^{-5}$. Scores of resulting BLASTp hits were normalized to the score when searched
158  against the set of predicted peptides. Only predicted sequences with at least 0.25 score ratio and 0.25
159  query length covered by the best alignment were kept. Representative transcript and peptide
160  sequences were identified per gene to encode the longest possible peptide as previously established
161  [29,35]. GO terms were assigned via InterProScan5 [36]. Reciprocal best BLAST hits against
162  Araport11 [35] were identified based on a previously developed script [27]. Remaining sequences
163  were annotated via best BLAST hits against Araport11 with an e-value cut-off at 0.0001. The
164  Araport11 annotation was transferred to predicted sequences.
165       Prediction of non-protein coding RNA genes like tRNA and rRNA genes was performed based
166  on tRNAscan-SE v2.0.3 [37,38] and INFERNAL (cmscan) v1.1.2 [39] based on Rfam13 [40].
167       RepeatModeler v2 [41] was deployed with default settings for the identification of repeat family
168  consensus sequences.
169

170  *2.4. Assembly and annotation assessment*

171       The percentage of phased and merged regions in the genome was assessed with the focus on
172  predicted genes. Based on Illumina and ONT read mappings, the average coverage depth per gene
173  was calculated. The distribution of these average values per gene allowed the classification of genes
174  as phased (haploid read depth) or merged (diploid read depth). As previous studies revealed that
175  Illumina short reads have a higher resolution for such coverage analysis [42], we focused on the
176  Illumina read data set for these analyses. Sequence variants were detected based on this read
177  mapping as previously described [43]. The number of heterozygous variants per gene was
178  calculated and compared between the groups of putatively phased and merged genes. Predicted
179  peptide sequences were compared against the annotation of other species including *A. thaliana* and
180  *D. rotundata* via OrthoFinder v2 [44].
181       Sequence reads and assembled sequences are available at ENA under the project ID ERP118030
182  (see File S1 for details). The assembly described in this manuscript is available under
183  GCA_902712375. Additional annotation files including the contigs assigned to organelle genomes
184  are available as a data publication from the institutional repository of Bielefeld University at
185  https://doi.org/10.4119/unibi/2941469.
186       Alleles covered by the fraction of phase-separated gene models were matched based on
187  reciprocal best BLAST hits of the coding sequences (CDSs) following a previously described
188  approach [27]. Alleles were considered a valid pair that represents a single gene if the second best
189  match displayed 99% or less of the score of the best match. A customized Python script for this allele
190  assignment is available on github (https://github.com/bpucker/yam).
191

192  **3. Results**

193     In total, we generated 66.7 Gbp of ONT reads data representing respectively about 218x
194 coverage of the estimated 322 Mbp haploid *D. dumetorum* genome. Read length N50 of the raw ONT
195 data set was 23 kbp and increased to 38 kbp through correction, trimming, and filtering.
196 Additionally, 13 Gbp of Illumina short read data (about 40x coverage) were generated. After all
197 polishing steps, the final assembly represents 485 Mbp of the highly heterozygous *D. dumetorum*
198 genome with an N50 of 3.2 Mbp (Table 1). Substantial improvement of the initial assembly through
199 various polishing steps was indicated by the increasing number of recovered BUSCOs (File S2). The
200 final assembly displayed more BUSCOs (92.30% out of 1440 included in the embryophyta data set,
201 see File S2 for details on the various BUSCO classes) compared to the publicly available genome
202 sequence assembly of *D. rotundata* (v0.1) for that we detected 81.70% BUSCOs with identical
203 parameters. Since there is no genetic map available for *D. dumetorum*, we transferred linkage group
204 assignments from *D. rotundata* to our assembly. In total, 206 contigs comprising 330 Mbp were
205 assigned to a linkage group, while 718 contigs remained unplaced with a total sequence of 155 Mbp
206 (File S3). One plastid and six mitochondrial contigs were identified based on sequence similarity to
207 *D. rotundata* organelle genome sequences (see https://doi.org/10.4119/unibi/2941469); the assignment
208 was confirmed by very high coverage in the read mapping. Our *D. dumetorum* plastid sequence
209 turned out to almost identical to the data recently provided for the *D. dumetorum* plastome [11].
210     Haploid genome size estimations based on k-mer distributions of the Illumina sequence reads
211 ranged from 215 Mbp (gce) over 254 Mbp (GenomeScope) to 350 Mbp (findGSE, MGSE) (File S4). The
212 differences between the estimates might be influenced by the repeat content of the *D. dumetorum*
213 genome (see below).
214
215

216

217 **Table 1**. Statistics of selected versions of the *D. dumetorum* genome assembly (see File S5 for a full table).

218

| | Initial assembly | Racon1 | Medaka2 | Pilon3 | Final |
|---|---|---|---|---|---|
| Number of contigs | 1,172 | 1,172 | 1,215 | 1,215 | 924 |
| Max. contig length [bp] | 20,187,448 | 20,424,333 | 17,910,017 | 17,878,854 | 17,878,854 |
| Assembly size [bp] | 501,985,705 | 508,061,170 | 507,215,754 | 506,184,192 | 485,115,345 |
| Assembly size without N [bp] | 501,985,705 | 508,061,170 | 507,215,754 | 506,184,192 | 485,115,345 |
| GC content | 37.74% | 37.66% | 37.87% | 37.59% | 37.57% |
| N50 [bp] | 3,896,882 | 3,930,287 | 2,598,889 | 2,593,751 | 3,190,870 |
| N90 [bp] | 136,614 | 138,199 | 137,206 | 136,754 | 156,407 |
| BUSCO (complete) | 85.70% | 89.80% | 91.90% | 92.30% | 92.30% |

219

220 Different gene prediction approaches were evaluated (File S6) leading to a final set of 35,269
221 protein-encoding gene models. The average gene model spans 4.3 kbp, comprises 6 exons and
222 encodes 455 amino acids (see File S6 for details). The gene prediction dataset for *D. dumetorum* is
223 further supported by the identification of 6,475 single copy orthologs between *D. dumetorum* and *D.
224 rotundata* as well as additional orthogroups (File S7). Based on these single copy orthologs, the
225 similarity of *D. dumetorum* and *D. rotundata* sequences was determined to be mostly above 80% (File
226 S8). If the phase separated allelic gene models were considered (Figure 1), 3,352 additional single
227 copy orthologs were detected. Functional annotation was assigned to 23,835 genes (File S9).
228 Additionally, 9,941 non-coding RNA gene models were predicted including 784 putative tRNA
229 genes (see https://doi.org/10.4119/unibi/2941469). Finally and in addition to gene models encoding
230 proteins and various RNA types, we identified 1,129 repeat consensus sequences with a combined
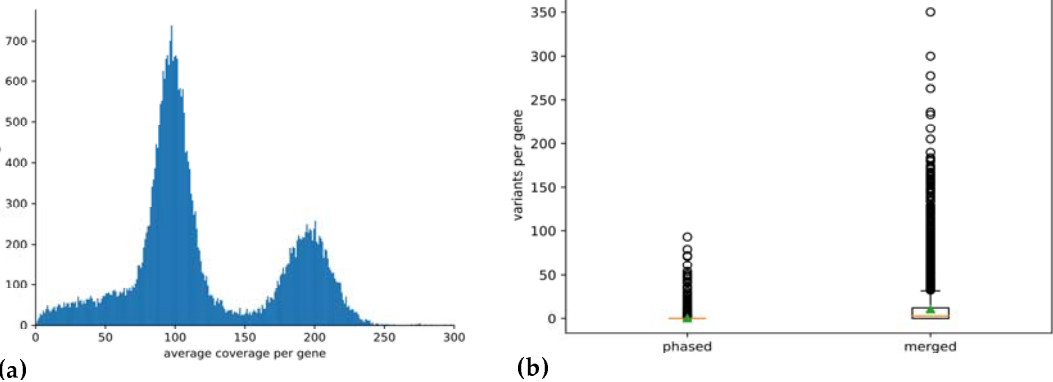231 length of 1.3 Mbp (File S10). The maximal repeat consensus length is 17.4 kbp, while the N50 is only
232 2.5 kbp.

233

234
235
236
237
238



(a)  (b)

239 **Figure 1**. (**a**) Distribution of the average sequencing read depth per gene models. Predicted gene models were
240 classified into phase separated and merged based on the average read depth value deduced from the analysis
241 presented here. The haploid read depth with Illumina short reads ranges from 50-fold to 150-fold. (**b**) Number
242 of heterozygous sequence variants in phase separated and merged genes. The high proportion of heterozygous

243 variants in merged gene models is due to the mapping of reads originating from two different alleles to the
244 same region of the assembly.

245

246   Average read mapping depth per gene was analyzed to distinguish genes annotated in
247 separated haplophases as well as merged sequences, respectively (Figure 1, File S11). About 64% of
248 all predicted protein-encoding gene models were in the expected range of the haploid read mapping
249 depth between 50-fold and 150-fold and about 27% are merged with a read depth between 150-fold
250 and 250-fold. Only 6% of all genes show an average read depth below 50-fold and only 1% show an
251 average coverage higher than 250-fold. It should be noted that the gene models annotated in the
252 phase separated part will cover in general two alleles per gene. A total of 22,885 gene models,
253 representing the 64% in the range of the haploid read mapping depth, were sorted into allelic pairs
254 which was successful for 8,492 genes. The findings presented above can be explained by a diploid
255 genome. An analysis with Smudgeplot indicated hints for a tetraploid genome from analysis with a
256 k-mer size of 19, while the other three investigated k-mer sizes supported a diploid genome (File
257 S12).

258

259

260 **4. Discussion**

261   The release of genome sequences of many model and crop plants has provided new
262 opportunities for gene identification and studies of genome evolution, both ultimately serving the
263 process of plant breeding [45] by allowing discovery of genes responsible for important agronomic
264 traits and the development of molecular markers associated with these traits. Here, we present the
265 first genome sequence for *Dioscorea dumetorum*, an important crop for Central and Western Africa,
266 and the second genome sequence for the genus. Our assembly offers a great opportunity to
267 understand the evolution of yam and to elucidate some biological constraints inherent to yam
268 including a long growth cycle, poor to non-flowering, polyploidy, vegetative propagation, and a
269 heterozygous genetic background [46]. Yam improvement has been challenging due to these factors
270 preventing the genetic study of important traits in yam [47].
271   Oxford Nanopore sequencing has proven to be a reliable and affordable technology for
272 sequencing genomes thus replacing Illumina technique for *de novo* genome sequencing due to
273 substantially higher assembly continuity [42,48]. Large fractions of the genome sequence were
274 separated into phases, while regions with lower heterozygosity are merged into one representative
275 sequence. Coverage analysis with Illumina read mapping allowed to classify predicted gene models
276 as 'phased' or 'merged' based on an average coverage around 100 fold or around 200 fold,
277 respectively. While this distinction is possible at the gene model level, whole contigs cannot be
278 classified this way. Several Mbp long contigs comprise alternating phase separated and merged
279 regions. Therefore, it is likely that the contigs represent a mixture of both haplophases with the risk
280 of switching between phases at each merged region. Since the haplophases cannot be resolved
281 continuously through low heterozygosity regions, purging of contigs to reduce the assembly into a
282 representation of the haploid genome might be advantageous for some applications in the future.
283 The bimodal coverage distribution (Figure 1a) supports the assumption that *D. dumetorum* Ibo sweet
284 3 has a diploid genome. This is supported by Smudgeplot for three out of 4 k-mer sizes tested while
285 the shortest k-mer size used (19) finds indications for tetraploidy. Since a high ploidy would result in
286 more distinct coverage peaks as observed for a genome with up to pentaploid parts [42], we assume
287 that the genome is diploid. The weak hint for tetraploidy might be due to a whole genome
288 duplication event early in the diversification of the genus. The N50 of 3.2 Mbp is in the expected
289 range for a long read assembly of a highly heterozygous plant species which contains quite some
290 repetitive sequences as others reported similar values before [49]. Due to regions of merged
291 haplophases the total assembly size of 485 Mbp is smaller than expected for a fully phase separated
292 "diploid" genome sequence based on the haploid genome size estimation of 322 Mbp.

293      We noticed an increase of the number of BUSCOs through several polishing rounds. Initial
294 assemblies of long reads can contain numerous short insertions and deletions as these are the major
295 error type in ONT reads [50]. As a result, the identification of CDSs and deduced open reading
296 frames is hindered through apparent disruptions of some CDS. Through the applied polishing steps,
297 the number of such apparent frame shifts is reduced thus leading to an increase of detected BUSCOs.
298      *D. dumetorum* has 36 chromosomes [12], so with 924 contigs we are far from chromosome-level
299 resolution but considerably better than the other genome assembly published in the genus, that of *D.*
300 *rotundata* with 40 chromosomes [8]. Knuth [51] circumscribed *D. dumetorum* and *D. rotundata* in two
301 distant sections *D.* sect. *Lasiophyton* and *D.* sect. *Enantiophlyllum*, respectively. Also, phylogenetically
302 the two species are quite distantly related with a last common ancestor about 30 million years ago
303 [11,52]. Comparing our predicted peptides to the *D. rotundata* peptide set [8], we identified about
304 9,800 single copy orthologs (6,475 in the whole set of 35,269 gene models plus 3,352 with a relation of
305 one gene in *D. rotundata* and two phase-separated alleles in *D. dumetorum*) which could elucidate the
306 evolutionary history of those species. The total number of predicted protein-encoding gene models
307 was determined to be 35,269, but this number includes two copies of about 11,300 gene models (see
308 Figure 1) as these are represented by two alleles each. The CDS-based pairing we performed
309 detected about 8,500 of the theoretical maximum of 11,300 cases which is a good success rate given
310 the fact that close paralogs and also hemizygous genome regions contribute to the detected number
311 of phase-separated gene models. If phase-separated gene models (alleles) are excluded, a number of
312 about 24,000 genes would result for *D. dumetorum*. This fits to the range detected in other higher
313 plant genomes [53,54]. The BUSCO results support this interpretation with about 40% of BUSCOs
314 that occur with exactly two copies. Therefore, the true number of protein-encoding genes of a
315 haploid *D. dumetorum* (trifoliate yam) genome could be around 25,000, also considering that the
316 BUSCO analysis indicated by 5.8% missing BUSCOs that still a small fraction of the genome
317 sequence is missing. This gene number fits well to gene numbers of higher plants based on all
318 available annotations at NCBI/EBI [54]. The average length of genes and the number of encoded
319 amino acids are in the same range as previously observed for other plant species from diverse
320 taxonomic groups [33,55].
321      It should be noted that the assignment of *D. dumetorum* sequences to the *D. rotundata*
322 pseudochromosomes and indirectly the respective linkage groups contains the risk of incorrect
323 assignments. However, although *D. rotundata* and *D. dumetorum* are evolutionary separated, *D.*
324 *rotundata* is the most closely related species with genetic and genomic resources.
325      Our draft genome has the potential to provide a complete new way to breed in *D. dumetorum*,
326 for example avoiding the post-harvest hardening phenomenon, which begins within 24 h after
327 harvest and makes it necessary to process the tubers within this time to allow consumption [2]. The
328 family Dioscoreaceae consists of more than 800 species [56] and the post-harvest hardening
329 phenomenon has only been reported from *D. dumetorum* [57], outlining the singularity of this species
330 among yam species. We predicted a large number of genes, which will include putative genes
331 controlling the post-harvest hardening on *D. dumetorum* and many useful bioactive compounds
332 detected in this yam species, which is considered the most nutritious and valuable from a
333 phytomedical point of view [58]. Ongoing work will try to identify these genes and polymorphisms
334 for making them available for subsequent breeding.
335      In summary, we present the first *de novo* nuclear genome sequence assembly of *D. dumetorum*
336 with very good contiguity and partially separated phases. Our assembly has no ambiguous bases
337 with a well applicable protein-encoding gene annotation. This assembly unraveled the genomic
338 structure of *D. dumetorum* to a large extent and will serve as a reference genome sequence for yam
339 breeding by helping to identify and develop molecular markers associated with relevant agronomic
340 traits, and to understand the evolutionary history of *D. dumetorum* and yam species in general.
341

342      **Supplementary Materials**: The following are available online:

343     File S1: Sequencing overview with ENA identifiers of runs.
344     File S2: Results of BUSCO analysis of different assembly versions.
345     File S3: AGP file describing contig assignment to *D. rotundata* pseudochromosomes.
346     File S4: Genome size estimation overview using four different tools.
347     File S5: General statistics of different assembly versions.
348     File S6: Comparison of results from different gene prediction approaches.
349     File S7: Orthogroups of predicted peptides of *D. rotundata* and *D. dumetorum*.
350     File S8: Similarity of *D. dumetorum* and *D. rotundata* based on single copy orthologs.
351     File S9: Functional annotation of predicted genes in the *D. dumetorum* genome sequence.
352     File S10: Consensus sequences of repeat elements detected in the *D. dumetorum* genome sequence.
353     File S11: Average short read mapping coverage of predicted genes in the *D. dumetorum* genome sequence.
354     File S12: Results from Smudgeplot analyses.
355

356     **Author Contributions**: CS, BP, DCA, and BW designed the study. CS collected the sample. BP performed DNA
357     extraction, ONT sequencing, and genome assembly. PV performed Illumina sequencing. CS and BP processed
358     the assembly. BP performed gene prediction and evaluation. CS and BP wrote the initial draft. BW and DCA
359     revised the manuscript. All authors read and approved the final version of the manuscript.

367

368     **Conflicts of Interest**: The authors declare no conflict of interest.

369

370

371     **References**

372

373     1.     Viruel, J.; Forest, F.; Paun, O.; Chase, M.W.; Devey, D.; Couto, R.S.; Segarra-Moragues, J.G.;
374            Catalan, P.; Wilkin, P. A nuclear Xdh phylogenetic analysis of yams (Dioscorea
375            Dioscoreaceae) congruent with plastid trees reveals a new Neotropical lineage. *Botanical*
376            *Journal of the Linnean Society* **2018**, 1-15.
377     2.     Sefa-Dedeh, S.; E.O., A. Biochemical and textural changes in trifoliate yam Dioscorea
378            dumetorum tubers after harvest. *Food Chemistry* **2002**, *79*, 27-40.
379     3.     Alozie, Y.E.; Akpanabiatu, M.; Eyong, E.U.; Umoh, I.B.; Alozie, G. Amino Acid Composition
380            of Dioscorea dumetorum Varities. *Pakistan Journal of Nutrition* **2009**, *8*, 103-105.
381     4.     Ferede, R.; Maziya-Dixon, B.; Alamu, O.E.; Asiedu, R. Identification and quantification of
382            major carotenoids of deep yellow-fleshed yam (tropical Dioscorea dumetorum). *Journal of*
383            *Food, Agruculture & Environment* **2010**, *8*, 160-166.
384     5.     Nimenibo-Uadia, R.; Oriakhi, A. Proximate, Mineral and Phytochemical Composition of
385            Dioscorea dumetorium Pax. *J. Appl. Sci. Environ. Manage.* **2017**, *21*, 771-774.
386     6.     Iwu, M.M.; Okunji, C.O.; Ohiaeri, G.O.; Akah, P.; Corley, D.; Tempesta, M.S. Hypoglycaemic
387            activity of dioscoretine from tubers of Dioscorea dumetorum in normal and alloxan diabetic
388            rabbits. *Planta Medica* **1990**, *56*, 264-267.

389 7.  Siadjeu, C.; Panyoo, E.A.; Toukam, G.M.S.; Bell, J.M.; Nono, B.; Medoua, G.N. Influence of
390     Cultivar on the Postharvest Hardening of Trifoliate Yam (Dioscorea dumetorum) Tubers.
391     *Hindawi* **2016**, *16*.
392 8.  Tamiru, M.; Natsume, S.; Takagi, H.; White, B.; Yaegashi, H.; Shimizu, M.; Yoshida, K.;
393     Uemura, A.; Oikawa, K.; Abe, A., et al. Genome sequencing of the staple food crop white
394     Guinea yam enables the development of a molecular marker for sex determination. *BMC*
395     *Biology* **2017**, *15*, 86.
396 9.  Cormier, F.; Lawac, F.; Maledon, E.; Gravillon, M.C.; Nudol, E.; Mournet, P.; Vignes, H.;
397     Chaïr, H.; Arnau, G. A reference high-density genetic map of greater yam (Dioscorea alata
398     L.). *Theoretical and Applied Genetics* **2019**, *132*, 1733-1744, doi:10.1007/s00122-019-03311-6.
399 10. Ngo Ngwe, M.F.; Omokolo, D.N.; Joly, S. Evolution and Phylogenetic Diversity of Yam
400     Species (Dioscorea spp.): Implication for Conservation and Agricultural Practices. *PLoS One*
401     **2015**, *10*, e0145364, doi:10.1371/journal.pone.0145364.
402 11. Magwe-Tindo, J.; Wieringa, J.J.; Sonke, B.; Zapfack, L.; Vigouroux, Y.; Couvreur, T.L.P.;
403     Scarcelli, N. Complete plastome sequences of 14 African yam species (Dioscorea spp.).
404     *Mitochondrial DNA Part B-Resources* **2019**, *4*, 74-76, doi:10.1080/23802359.2018.1536466.
405 12. Miege, J. Nombres chromosomiques et répartition géographique de quelques plantes
406     tropicales et équatoriales. *Revue de Cytologie et de Biologie Végétales* **1954**, *15*, 312-348.
407 13. Hui-Chen, C.; Mei-Chen, C.; Ping-Ping, L.; Chih-Tsun, T.; Fang-Ping, D. A cytotaxonomic
408     study on Chinese Dioscorea L. - the chromosome numbers and their relation to the origin
409     and evolution of the genus. *Journal of Systematics and Evolution* **1985**, *23*, 11-18.
410 14. Siadjeu, C.; Mayland-Quellhorst, E.; Albach, D.C. Genetic diversity and population structure
411     of trifoliate yam (Dioscorea dumetorum Kunth) in Cameroon revealed by
412     genotyping-by-sequencing (GBS). *BMC Plant Biology* **2018**, *18*, 359.
413 15. Rosso, M.G.; Li, Y.; Strizhov, N.; Reiss, B.; Dekker, K.; Weisshaar, B. An *Arabidopsis thaliana*
414     T-DNA mutagenised population (GABI-Kat) for flanking sequence tag based reverse
415     genetics. *Plant Molecular Biology* **2003**, *53*, 247-259.
416 16. Vurture, G.W.; Sedlazeck, F.J.; Nattestad, M.; Underwood, C.J.; Fang, H.; Gurtowski, J.;
417     Schatz, M.C. GenomeScope: fast reference-free genome profiling from short reads.
418     *Bioinformatics* **2017**, *33*, 2202-2204, doi:10.1093/bioinformatics/btx153.
419 17. Sun, H.; Ding, J.; Piednoel, M.; Schneeberger, K. findGSE: estimating genome size variation
420     within human and Arabidopsis using k-mer frequencies. *Bioinformatics* **2018**, *34*, 550-557,
421     doi:10.1093/bioinformatics/btx637.
422 18. Liu, B.; Shi, Y.; Yuan, J.; Hu, X.; Zhang, H.; Li, N.; Li, Z.; Chen, Y.; Mu, D.; Fan, W. Estimation
423     of genomic characteristics by analyzing k-mer frequency in de novo genome projects. *arXiv*
424     **2013**, 1308.2012 [q-bio.GN].
425 19. Marcais, G.; Kingsford, C. A fast, lock-free approach for efficient parallel counting of
426     occurrences of k-mers. *Bioinformatics* **2011**, *27*, 764-770, doi:10.1093/bioinformatics/btr011.
427 20. Pucker, B. Mapping-based genome size estimation. *bioRxiv* **2019**, 10.1101/607390,
428     doi:10.1101/607390.
429 21. Ranallo-Benavidez, T.R.; Jaron, K.S.; Schatz, M.C. GenomeScope 2.0 and Smudgeplots:
430     Reference-free profiling of polyploid genomes. *bioRxiv* **2019**, 10.1101/747568,
431     doi:10.1101/747568.

432   22.   Koren, S.; Walenz, B.P.; Berlin, K.; Miller, J.R.; Bergman, N.H.; Phillippy, A.M. Canu:
433         scalable and accurate long-read assembly via adaptive k-mer weighting and repeat
434         separation. *Genome Research* **2017**, *27*, 722-736.
435   23.   Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **2018**, *34*,
436         3094-3100.
437   24.   Vaser, R.; Sović, I.; Nagarajan, N.; Šikić, M. Fast and accurate de novo genome assembly
438         from long uncorrected reads. *Genome Research* **2017**, *27*, 737-746.
439   25.   Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
440         *arXiv* **2013**, 1303.3997v1302 (Preprint posted May 1326, 2013).
441   26.   Walker, B.J.; Abeel, T.; Shea, T.; Priest, M.; Abouelliel, A.; Sakthikumar, S.; Cuomo, C.A.;
442         Zeng, Q.; Wortman, J.; Young, S.K., et al. Pilon: an integrated tool for comprehensive
443         microbial variant detection and genome assembly improvement. *PLoS ONE* **2014**, *9*, e112963.
444   27.   Pucker, B.; Holtgräwe, D.; Rosleff Sörensen, T.; Stracke, R.; Viehöver, P.; Weisshaar, B. A De
445         Novo Genome Sequence Assembly of the Arabidopsis thaliana Accession Niederzenz-1
446         Displays Presence/Absence Variation and Strong Synteny. *PLoS ONE* **2016**, *11*, e0164321.
447   28.   Tang, H.; Zhang, X.; Miao, C.; Zhang, J.; Ming, R.; Schnable, J.C.; Schnable, P.S.; Lyons, E.;
448         Lu, J. ALLMAPS: robust scaffold ordering based on multiple maps. *Genome Biol* **2015**, *16*, 3,
449         doi:10.1186/s13059-014-0573-1.
450   29.   Pucker, B.; Holtgräwe, D.; Weisshaar, B. Consideration of non-canonical splice sites
451         improves gene prediction on the Arabidopsis thaliana Niederzenz-1 genome sequence. *BMC*
452         *Research Notes* **2017**, *10*, 667.
453   30.   Simão, F.A.; Waterhouse, R.M.; Ioannidis, P.; Kriventseva, E.V.; Zdobnov, E.M. BUSCO:
454         assessing genome assembly and annotation completeness with single-copy orthologs.
455         *Bioinformatics* **2015**, *31*, 3210-3212.
456   31.   Keller, O.; Kollmar, M.; Stanke, M.; Waack, S. A novel hybrid gene prediction method
457         employing protein multiple sequence alignments. *Bioinformatics* **2011**, *27*, 757-763.
458   32.   Smit, A.F.A.; Hubley, R.; Green, P. RepeatMasker Open-4.0. Availabe online:
459         http://www.repeatmasker.org (accessed on
460   33.   Pucker, B.; Feng, T.; Brockhington, S. Next generation sequencing to investigate genomic
461         diversity in Caryophyllales. *bioRxiv* **2019**, , doi:10.1101/646133 (Preprint posted 2019-07-27).
462   34.   Altschul, S.F.; Gish, W.; Miller, W.; Myers, E.W.; Lipman, D.J. Basic local alignment search
463         tool. *Journal of Molecular Biology* **1990**, *215*, 403-410.
464   35.   Cheng, C.Y.; Krishnakumar, V.; Chan, A.; Thibaud-Nissen, F.; Schobel, S.; Town, C.D.
465         Araport11: a complete reannotation of the Arabidopsis thaliana reference genome. *The Plant*
466         *Journal* **2017**, 789-804.
467   36.   Jones, P.; Binns, D.; Chang, H.Y.; Fraser, M.; Li, W.; McAnulla, C.; McWilliam, H.; Maslen, J.;
468         Mitchell, A.; Nuka, G., et al. InterProScan 5: genome-scale protein function classification.
469         *Bioinformatics* **2014**, *30*, 1236-1240, doi:10.1093/bioinformatics/btu031.
470   37.   Lowe, T.M.; Eddy, S.R. tRNAscan-SE: a program for improved detection of transfer RNA
471         genes in genomic sequence. *Nucleic Acids Research* **1997**, *25*, 955-964.
472   38.   Chan, P.P.; Lowe, T.M. tRNAscan-SE: Searching for tRNA genes in genomic sequences. In
473         *Gene Prediction: Methods and Protocols*, 2019/04/26 ed.; Kollmar, M., Ed. Springer New York:
474         New York, 2019; Vol. 1962, pp. 1-14.

475    39.    Nawrocki, E.P.; Eddy, S.R. Infernal 1.1: 100-fold faster RNA homology searches.
476          *Bioinformatics* **2013**, *29*, 2933-2935.

477    40.    Kalvari, I.; Argasinska, J.; Quinones-Olvera, N.; Nawrocki, E.P.; Rivas, E.; Eddy, S.R.;
478          Bateman, A.; Finn, R.D.; Petrov, A.I. Rfam 13.0: shifting to a genome-centric resource for
479          non-coding RNA families. *Nucleic Acids Research* **2018**, *46*, D335-D342.

480    41.    Flynn, J.M.; Hubley, R.; Goubert, C.; Rosen, J.; Clark, A.G.; Feschotte, C.; Smit, A.F.
481          RepeatModeler2: automated genomic discovery of transposable element families. *bioRxiv*
482          **2019**, 10.1101/856591, doi:10.1101/856591.

483    42.    Pucker, B.; Ruckert, C.; Stracke, R.; Viehover, P.; Kalinowski, J.; Weisshaar, B. Twenty-Five
484          Years of Propagation in Suspension Cell Culture Results in Substantial Alterations of the
485          Arabidopsis Thaliana Genome. *Genes* **2019**, *10*, 671, doi:10.3390/genes10090671.

486    43.    Baasner, J.S.; Howard, D.; Pucker, B. Influence of neighboring small sequence variants on
487          functional impact prediction. *bioRxiv* **2019**, , doi:10.1101/596718 (Preprint posted 2019-06-13).

488    44.    Emms, D.M.; Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
489          genomics. *Genome Biology* **2019**, *20*, 238, doi:10.1186/s13059-019-1832-y.

490    45.    Ruggieri, V.; Alexiou, K.G.; Morata, J.; Argyris, J.; Pujol, M.; Yano, R.; Nonaka, S.; Ezura, H.;
491          Latrasse, D.; Boualem, A., et al. An improved assembly and annotation of the melon
492          (Cucumis melo L.) reference genome. *Scientic Reports* **2018**, *8*, 8088,
493          doi:10.1038/s41598-018-26416-2.

494    46.    Mignouna, H.D.; Abang, M.M.; Asiedu, R. Harnessing modern biotechnology for tropical
495          tuber crop improvement: Yam (Dioscorea spp.) molecular breeding. *African Journal of*
496          *Biotechnology* **2003**, *2*, 12.

497    47.    Mignouna, H.D.; Abang, M.M.; Asiedu, R. Genomics of Yams, a Common Source of Food
498          and Medicine in the Tropics. In *Genomics of Tropical Crop Plants*, Moore, P.H., Ming, R., Eds.
499          2008; 10.1007/978-0-387-71219-2_23pp. 549-570.

500    48.    Michael, T.P.; Jupe, F.; Bemm, F.; Motley, S.T.; Sandoval, J.P.; Lanz, C.; Loudet, O.; Weigel,
501          D.; Ecker, J.R. High contiguity Arabidopsis thaliana genome assembly with a single
502          nanopore flow cell. *Nature Communications* **2018**, *9*, 541.

503    49.    Paajanen, P.; Kettleborough, G.; Lopez-Girona, E.; Giolai, M.; Heavens, D.; Baker, D.; Lister,
504          A.; Cugliandolo, F.; Wilde, G.; Hein, I., et al. A critical comparison of technologies for a plant
505          genome sequencing project. *Gigascience* **2019**, *8*, doi:10.1093/gigascience/giy163.

506    50.    Salmela, L.; Walve, R.; Rivals, E.; Ukkonen, E. Accurate self-correction of errors in long reads
507          using de Bruijn graphs. *Bioinformatics* **2017**, *33*, 799-806, doi:10.1093/bioinformatics/btw321.

508    51.    Knuth, R. Dioscoreaceae. In *Das Pflanzenreich*, Engelr, A., Ed. Engelmann, W.: Leipzig, 1924.

509    52.    Viruel, J.; Segarra-Moragues, J.G.; Raz, L.; Forest, F.; Wilkin, P.; Sanmartin, I.; Catalan, P.
510          Late Cretaceous-Early Eocene origin of yams (Dioscorea, Dioscoreaceae) in the Laurasian
511          Palaearctic and their subsequent Oligocene-Miocene diversification. *Journal of Biogeography*
512          **2016**, *43*, 750-762, doi:10.1111/jbi.12678.

513    53.    Wendel, J.F.; Jackson, S.A.; Meyers, B.C.; Wing, R.A. Evolution of plant genome architecture.
514          *Genome Biology* **2016**, *17*, 37, doi:10.1186/s13059-016-0908-1.

515    54.    Pucker, B.; Brockington, S.F. Genome-wide analyses supported by RNA-Seq reveal
516          non-canonical splice sites in plant genomes. *BMC Genomics* **2018**, *19*, 980.

517    55.    Pucker, B.; Holtgräwe, D.; Stadermann, K.B.; Frey, K.; Huettel, B.; Reinhardt, R.; Weisshaar,
518           B. A chromosome-level sequence assembly reveals the structure of the Arabidopsis thaliana
519           Nd-1 genome and its gene set. *PLoS One* **2019**, *14*, e0216233.
520    56.    Barton, H. Yams: Origins and Development. **2014**.
521    57.    Treche, S.; Delpeuch, F. Physiologie Vegetale - Mise en evidence de l'apparition d'un
522           epaississement membranaire dans le parenchyme des tubercules de Dioscorea dumetorum
523           au cours de la conservation. *C. R. Acad. Sc. Paris* **1979**, *288*.
524    58.    Price, E.J.; Wilkin, P.; Sarasan, V.; Fraser, P.D. Metabolite profiling of Dioscorea (yam)
525           species reveals underutilised biodiversity and renewable sources for high-value
526           compounds. *Scientic Reports* **2016**, *6*.

527