# High contiguity *de novo* genome sequence assembly of Trifoliate yam (*Dioscorea dumetorum*) using long read sequencing

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

[1] Institute for Biology and Environmental Sciences, Biodiversity and Evolution of Plants, Carl-von-Ossietzky University Oldenburg, Carl-von-Ossietzky Str. 9-11, 26111 Oldenburg, Germany; christian.siadjeu@uol.de (CS); dirk.albach@uol.de (DCA)

[2] Genetics and Genomics of Plants, Faculty of Biology, Center for Biotechnology (CeBiTec), Bielefeld University, Sequenz 1, 33615 Bielefeld, NRW, Germany; bpucker@cebitec.uni-bielefeld.de (BP), viehoeve@cebitec.uni-bielefeld.de (PV), bernd.weisshaar@uni-bielefeld.de (BW)

[3] Molecular Genetics and Physiology of Plants, Faculty of Biology and Biotechnology, Ruhr-University Bochum, Universitätsstraße 150, 44801 Bochum, Germany; bpucker@cebitec.uni-bielefeld.de (BP)

* Correspondence: bernd.weisshaar@uni-bielefeld.de; Tel.: +49-521-106-8720

+ Shared first authorship. CS and BP contributed equally to this work and are co-first authors.

**Abstract:** The yam species *Dioscorea dumetorum* is one example of an orphan crop, not traded internationally. Post-harvest hardening starts within 24 hours after harvesting and renders the tubers inedible. Genomic resources are required for trifoliate yam to improve breeding for non-hardening varieties and for other traits. Here, we describe the sequencing of the *D. dumetorum* genome and the generation of a *de novo* assembly together with a corresponding annotation. The two haplophases of this highly heterozygous genome are separated to a large extent. The assembly represents 485 Mbp of the genome with an N50 of over 3.2 Mbp. A total of 35,269 protein-encoding gene structures as well as 9,941 non-coding RNA genes were predicted and functional annotations were assigned.

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

## 1. Introduction

The yam species *Dioscorea dumetorum* (trifoliate yam) belongs to the genus *Dioscorea* comprising about 600 described species. The genus is widely distributed throughout the tropics [1] and includes important root crops that offer staple food for over 300 million people. Eight *Dioscorea* species are commonly consumed in West and Central Africa, of which *D. dumetorum* has the highest nutrient value [2]. Yet, *D. dumetorum* constitutes an underutilized and neglected species despite its great potential for nutritional, agricultural, and pharmaceutical purposes. Tubers of *D. dumetorum* are protein-rich (9.6%) with a fairly balanced essential amino acids composition [3]. The provitamin A and carotenoids contents of the tubers of deep yellow genotypes are equivalent to those of yellow corn maize lines selected for increased concentrations of provitamin A [4]. The deep yellow yam tubers are used in antidiabetic treatments in Nigeria [5], probably due to the presence of dioscoretine, which is a bioactive compound with hypoglycaemic properties [6].

Unlike other yam species, the cultivation of *D. dumetorum* is limited by post-harvest hardening, which starts within 24 h after harvest and renders tubers inedible. Previous research showed that among 32 *D. dumetorum* cultivars tested, one cultivar was not affected by the hardening

44  phenomenon [7]. This discovery provides a starting point for a breeding program of *D. dumetorum*
45  against the post-harvest hardening phenomenon. *Dioscorea* cultivars are obligate outcrossing plants
46  that display highly heterozygous genomes. Thus, methods of genetic analysis routinely used in
47  inbreeding species such as linkage analysis using the segregation progeny of an F2 generation and
48  recombinant inbred lines are inapplicable to yam [8]. Furthermore, the development of
49  marker-assisted selection requires the establishment of marker assays and dense genetic linkage
50  maps. Thus, access to a complete and well-annotated genome sequence is one essential step towards
51  the implementation of up-to-date genetic and genomic approaches for *D. dumetorum* breeding.
52  Although the plastome sequences of 14 *Dioscorea* species [9] and the nuclear genome sequence of
53  *Dioscorea rotundata* [8] have been published, there is still a need for data from additional,
54  phylogenetically unrelated yam species. Here, we report sequencing and *de novo* assembly of the *D.*
55  *dumetorum* Ibo sweet 3 genome sequence based on long reads.
56

57  **2. Materials and Methods**

58  *2.1. Sampling and Sequencing*

59  A diploid *D. dumetorum* accession Ibo sweet 3 that does not display post-harvest hardening was
60  collected in the South-West region of Cameroon in 2013 [7]. Tubers of this accession were transferred
61  to Oldenburg (Germany) and the corresponding plants were cultivated in a greenhouse at 25°C. The
62  haploid genome size of the Ibo sweet 3 genotype had been estimated to be 322 Mbp through flow
63  cytometry [10].
64  DNA was extracted from 1g of leaf tissue using a CTAB-based method modified from [11].
65  After grinding the sample in liquid nitrogen, the powder was suspended in 5mL CTAB1 buffer
66  supplemented with 300μL ß-mercaptoethanol. The suspension was incubated at 75°C for 30 minutes
67  and inverted every five minutes. Next, 5mL dichloromethane were added and the solutions were
68  mixed by inverting. The sample was centrifuged at 11,200 g at 20°C for 30 minutes. The clear
69  supernatant was mixed with 10mL CTAB2 in a new reaction tube by inverting. Next, a
70  centrifugation was performed at 11,200 g at 20°C for 30 minutes. After discarding the supernatant,
71  1mL NaCl (1M) was added to re-suspend the sediment by gently flicking the tube. By adding an
72  equivalent amount of 1mL isopropanol and careful mixing, the DNA was precipitated again and the
73  sample was centrifuged as described above. After washing the sediment with 1mL of 70% ethanol,
74  200μL CTAB-TE buffer containing 2 mg RNaseA were added. Re-suspension and RNA degradation
75  were achieved by incubation over night at room temperature. DNA quality and quantity were
76  assessed via NanoDrop2000 measurement, agarose gel electrophoresis, and Qubit measurement.
77  The SRE kit (Circulomics) was used to enrich long DNA fragments following the suppliers'
78  instructions. Results were validated via Qubit measurement.
79  Library preparation was performed with 1μg of high molecular weight DNA following the
80  SQK-LSK109 protocol (Oxford Nanopore Technologies, ONT). Sequencing was performed on four
81  R9.4.1 flow cells on a GridION. Flow cells were treated with nuclease flush (20μL DNaseI (NEB) and
82  380μL nuclease flush buffer), once the number of active pores dropped below 200, to allow
83  successive sequencing of multiple libraries on an individual flow cell. Live base calling was
84  performed on the GridION by Guppy (ONT).
85  The Illumina paired-end (PE) library preparation was performed according to the Illumina
86  TruSeq DNA Sample Preparation v2 Guide. High molecular weight DNA was fragmented by
87  nebulization. End pair and A-tailing adaptors were ligated to the fragmented DNA. A two percent
88  low melt agarose gel was used to size select adaptor-ligated fragments. The fragments harbouring
89  adaptors on the both ends were enriched by PCR and final libraries were evaluated using PicoGreen.
90  Average fragment size of the libraries was estimated on a Bio Analyzer High Sensitivity DNA chip.
91  The PE library with an insert size of 700 to 790 bp was sequenced with 2 x 250 nt mode on an
92  Illumina HiSeq-1500.
93

94 *2.2. Genome assembly and polishing*

95     Canu v1.8 [12] was deployed for the genome assembly with the following parameters
96 "'genomeSize = 350m', 'corOutCoverage = 200' 'correctedErrorRate = 0.12' batOptions = -dg 3 -db 3
97 -dr 1 -ca 500 -cp 50' 'minReadLength = 10000' 'minOverlapLength = 5000' 'corMhapFilterThreshold =
98 0.0000000002' 'ovlMerThreshold = 500' 'corMhapOptions = --threshold 0.85 –num-hashes 512
99 –num-min-matches 3 –ordered-sketch-size 1000 –ordered-kmer-size 14 –min-olap-length 5000
100 –repeat-idf-scale 50'". The value for the genome size, estimated to be 322 Mbp, was increased to 350
101 Mbp to increase the number of reads utilized for the assembly process.
102     ONT reads were mapped back to the assembled sequence with minimap2 v2.17 [13], using the
103 settings recommended for ONT reads. Next, the contigs were polished by racon v.1.4.7 [14] with -m
104 8 -x -6 -g -8 as recommended prior to the polishing step with medaka. Two runs of medaka v.0.10.0
105 (https://github.com/nanoporetech/medaka) polishing were performed with default parameters (-m
106 r941_min_high) using ONT reads. Illumina short reads were aligned to the medaka consensus
107 sequence using BWA-MEM v. 0.7.17 [15]. This alignment was subjected to Pilon v1.23 [16] for final
108 polishing in three iterative rounds.
109     Downstream processing was based on a previously described workflow [17] and performed by
110 customized Python scripts for purging of short contigs (<100kb)  and calculation of assembly
111 statistics (https://github.com/bpucker/yam). Contigs with less than 3-fold average coverage in an
112 Illumina short read mapping were compared against nt via BLASTn with an e-value cut-off at $10^{-10}$ to
113 identify and remove bacterial and fungal sequences.
114

115 *2.3. Genome sequence annotation*

116     Hints for gene prediction were generated by aligning *D. rotundata* transcript sequences (TDr96
117 v1.0 [8] as previously described [18]. BUSCO v3 [19] was applied to generate a species-specific
118 parameter set. For comparison, the *D. rotundata* genome assembly GCA_002260605.1 [8] was
119 retrieved from NCBI. Hints and parameters were subjected to AUGUSTUS v.3.3 [20] for gene
120 prediction with previously described parameters [18]. Various approaches involving parameter files
121 of rice and maize as well as running the gene prediction on a sequence with masked repeats were
122 evaluated. BUSCO was applied again to assess the completeness of the gene prediction. The best
123 results for *D. dumetorum* genome sequence annotation were obtained based on an unmasked
124 assembly sequence with yam specific parameters generated via BUSCO as previously described
125 [19,21]. Predicted genes were filtered based on sequence similarity to entries in several databases
126 (UniProt/SwissProt, Araport11, *Brachypodium distachyon* v3.0, *Elaeis guineensis* v5.1,
127 GCF_000005425.2, GCF_000413155.1, *Musa acuminata* Pahang v2). Predicted peptide sequences were
128 compared to these databases via BLASTp [22] using an e-value cut-off of $10^{-5}$. Scores of resulting
129 BLASTp hits were normalized to the score when searched against the set of predicted peptides. Only
130 predicted sequences with at least 0.25 score ratio and 0.25 query length covered by the best
131 alignment were kept. Functional annotation was assigned via InterProScan5 [23] and through
132 sequence similarity to well characterized sequences. Representative transcript and peptide
133 sequences were identified per gene based on previously defined criteria to encode the longest
134 possible peptide [24, Pucker, 2017 #5337].
135     Prediction of non-protein coding RNA genes like tRNA and rRNA genes was performed based
136 on tRNAscan-SE v2.0.3 [25,26] and INFERNAL (cmscan) v1.1.2 [27] based on the Rfam13 [28].

137

138 *2.4. Assembly and annotation assessment*

139     The percentage of phased and merged regions in the genome was assessed with the focus on
140 predicted genes. Based on Illumina and ONT read mappings, the average coverage depth per gene
141 was calculated. The distribution of these average values per gene allowed the classification of genes

142 as phased (haploid read depth) or merged (diploid read depth). As previous studies revealed that
143 Illumina short reads have a higher resolution for such coverage analysis [29], we focused on the
144 Illumina read data set for these analyses. Sequence variants were detected based on this read
145 mapping as previously described [30]. The number of heterozygous variants per gene was
146 calculated and compared between the groups of putatively phased and merged genes. Predicted
147 peptide sequences were compared against the annotation of other species including *A. thaliana* and
148 *D. rotundata* via OrthoFinder v2 [31].

149 Sequence reads and assembled sequences are available at ENA under the project ID ERP118030
150 (see File S1 for details). The assembly described in this manuscript is available under
151 GCA_902712375. Additional annotation files are available from
152 https://docs.cebitec.uni-bielefeld.de/s/ArHmB4J2MXMsA5S.

153 Alleles covered by the fraction of phase separated gene structures were matched based on
154 reciprocal best BLAST hits of the coding sequences following a previously described approach [17].
155 Alleles were considered a valid pair that represents a single gene if the second best match displayed
156 99% or less of the score of the best match. A customized Python script for this allele assignment is
157 available on github (https://github.com/bpucker/yam).

158
159

## 3. Results

161 In total, we generated 70 Gbp of ONT reads data representing respectively about 218x coverage
162 of the estimated 322 haploid genome. Additionally, 13 Gbp of Illumina short read data (about 40x
163 coverage) were generated. After all polishing steps, the final assembly represents 485 Mbp of the
164 highly heterozygous *D. dumetorum* genome with an N50 of 3.2 Mbp (Table 1). Substantial
165 improvement of the initial assembly through various polishing steps is indicated by the increasing
166 number of recovered BUSCOs (File S2). The final assembly displayed more BUSCOs (92.30% out of
167 1440 included in the embryophyta data set) compared to the publicly available genome sequence
168 assembly of *D. rotundata* (v0.1) for that we detected 81.70% BUSCOs with identical parameters.

169

170 **Table 1**. Statistics of selected versions of the *D. dumetorum* genome assembly (see File S3 for a full table).

171

| | 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% |

172

173 Different gene prediction approaches were evaluated (File S4) leading to a final set of 35,269
174 protein-encoding gene structures. The average gene spans 4.3 kbp, comprises 6 exons and encodes
175 455 amino acids (see File S4 for details). The gene prediction dataset for *D. dumetorum* is further
176 supported by the identification of 6,475 single copy orthologs between *D. dumetorum* and *D.*
177 *rotundata* as well as additional orthogroups (File S5). If the phase separated allelic gene structures

178 were considered (Figure 1), 3,352 additional single copy orthologs were detected. Functional
179 annotation was assigned to 23,835 genes (File S6). Additionally, 9,941 non-coding RNA genes were
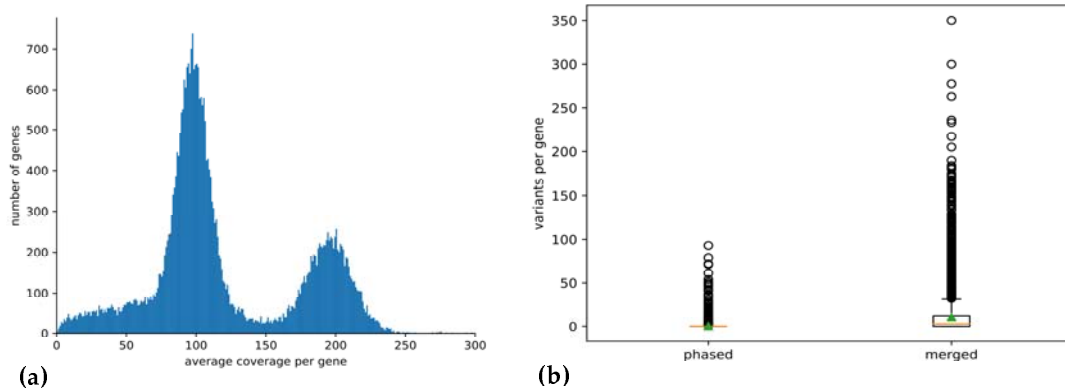180 predicted including 784 putative tRNA genes
181 (https://docs.cebitec.uni-bielefeld.de/s/ArHmB4J2MXMsA5S).

182

183



(a)                                                                (b)

188 **Figure 1**. (**a**) Distribution of the average sequencing read depth per gene structure. Predicted gene structures
189 were classified into phase separated and merged based on the average read depth value deduced from the
190 analysis presented here. The haploid read depth with Illumina short reads ranges from 50-fold to 150-fold. (**b**)
191 Number of heterozygous sequence variants in phase separated and merged genes. The high proportion of
192 heterozygous variants in merged genes is due to the mapping of reads originating from two different alleles to
193 the same region of the assembly.

194

195

196 Average read mapping depth per gene was analyzed to distinguish genes annotated in
197 separated and merged haplophases, respectively (Figure 1, File S7). About 64% of all predicted
198 protein-encoding gene structures were in the expected range of the haploid read mapping depth
199 between 50-fold and 150-fold and about 27% are merged with a read depth between 150-fold and
200 250-fold. Only 6% of all genes show an average read depth below 50-fold and only 1% show an
201 average coverage higher than 250-fold. It should be noted that the gene structures annotated in the
202 phase separated part will cover in general two alleles per gene. A total of 22,885 gene structures,
203 representing the 64% in the range of the haploid read mapping depth, were sorted into allelic pairs
204 which was successful for 8,492 genes.

205

206

207 **4. Discussion**

208 The release of genome sequences of many model and crop plants has provided new
209 opportunities for gene identification and studies of genome evolution, both ultimately serving the
210 process of plant breeding [32] by allowing discovery of genes responsible for important agronomic
211 traits and the development of molecular markers associated with these traits. Here, we present the
212 first genome sequence for *Dioscorea dumetorum*, an important crop for Central and Western Africa,
213 and the second for the genus. Our assembly offers a great opportunity to understand the evolution
214 of yam and to elucidate some biological constraints inherent to yam including a long growth cycle,
215 poor to non-flowering, polyploidy, vegetative propagation, and a heterozygous genetic background
216 [33]. Yam improvement has been challenging due to these factors preventing the genetic study of
217 important traits in yam [34].

218 Oxford Nanopore sequencing has proven to be a reliable and affordable technology for
219 sequencing genomes thus replacing Illumina technique for *de novo* genome sequencing due to
220 substantially higher assembly continuity [29,35]. Large fractions of the genome sequence were
221 separated into phases, while regions with lower heterozygosity are merged into one representative
222 sequence. Coverage analysis with Illumina read mapping allowed to classify predicted gene
223 structures as 'phased' or 'merged' based on an average coverage around 100 fold or around 200 fold,
224 respectively. While this distinction is possible at the gene structure level, whole contigs cannot be
225 classified this way. Several million bp long contigs comprise alternating phase separated and
226 merged regions. Therefore, it is likely that the contigs represent a mixture of both haplophases with
227 the risk of switching between phases at each merged region. Since the haplophases cannot be
228 resolved continuously through low heterozygosity regions, purging of contigs to reduce the
229 assembly into a representation of the haploid genome might be advantageous for some applications
230 in the future. The bimodal coverage distribution (Figure 1a) supports the assumption that *D.*
231 *dumetorum* Ibo sweet 3 has a diploid genome. A higher ploidy would result in more distinct coverage
232 peaks as observed for a genome with up to pentaploid parts [29]. The N50 of 3.2 Mbp is in the
233 expected range for a long read assembly of a highly heterozygous species as others reported similar
234 values before [36]. Due to regions of merged haplophases the total assembly size of 485 Mbp is
235 smaller than expected for a fully phase separated "diploid" genome sequence based on the haploid
236 genome size estimation of 322 Mbp.
237 Interestingly, we noticed an increase of the number of BUSCOs through several polishing
238 rounds. Initial assemblies of long reads can contain numerous short insertions and deletions as these
239 are the major error type [37]. The identification of open reading frames is hindered through apparent
240 disruptions of open reading frames. Through the applied polishing steps, the number of such
241 apparent frame shifts is reduced thus leading to an increase of detected BUSCOs.
242 *Dioscorea dumetorum* has 36 chromosomes [38], so with 924 contigs we are far from
243 chromosome-level resolution but considerably better than the other genome assembly published in
244 the genus, that of *D. rotundata* with 40 chromosomes [8]. Knuth [39], circumscribed *D. dumetorum*
245 and *D. rotundata* in two sections *Lasiophyton* and *Enantiophlyllum* respectively. Also, phylogenetically
246 the two species are quite distantly related [9]. Comparing our predicted peptides to the *D. rotundata*
247 peptide set [8], we identified about 9,800 single copy orthologs (6,475 in the whole set of 35,269 gene
248 structures plus 3,352 with a relation of one gene in *D. rotundata* and two phase separated alleles in *D.*
249 *dumetorum*) which could elucidate the evolutionary history of those species. Our genome sequence is
250 structurally accurate and more protein-encoding genes were predicted. The number of predicted
251 protein-encoding gene structures was determined to be 35,269, but this number includes two times
252 about 11,300 genes (see Figure 1) represented by two alleles. The CDS-based pairing we performed
253 detected about 8,500 of the theoretical maximum of 11.300 cases which is a good success rate given
254 the fact that close paralogs and also hemizygous genome regions contribute to the detected number
255 of phase separated gene structures. If phase separated gene structures (alleles) are excluded, a
256 number of about 24,000 genes would result for *D. dumetorum*. This fits to the range detected in other
257 higher plant genomes [40, Pucker, 2019 #5484]. The BUSCO results support this interpretation with
258 about 40% of BUSCOs that occur with exactly two copies. Therefore, the true number of
259 protein-encoding genes of a haploid yam genome could be around 25.000, also considering that the
260 BUSCO analysis indicated that still a small fraction of the genome sequence is missing. This gene
261 number fits well to gene numbers of higher plants based on all available annotations at the NCBI
262 [41] and is larger than that of *Theobroma cacao, Jatropha curcas, Oryza brchyantha*, and *Ananas comosus*.
263 The average length of genes and the number of encoded amino acids are in the same range as
264 previously observed for other plant species from diverse taxonomic groups [21,42].
265 Our draft genome has the potential to provide a complete new way to breed in *D. dumetorum*,
266 for example avoiding the post-harvest hardening phenomenon, which begins within 24 h after
267 harvest and makes it necessary to process the tubers within this time to allow consumption [2]. The
268 family Dioscoreaceae consists of more than 800 species [43] and the post-harvest hardening
269 phenomenon has only been reported from *D. dumetorum* [44], outlining the singularity of this species
270 among yam species. We predicted a large number of genes, which will include putative genes

271 controlling the post-harvest hardening on *D. dumetorum* and many useful bioactive compounds
272 detected in this yam species, which is considered the most nutritious and valuable from a
273 phytomedical point of view [45]. Ongoing work will try to identify these genes and polymorphisms
274 for making them available for subsequent breeding.

275     In summary, we present the first *de novo* nuclear genome sequence assembly of *D. dumetorum*
276 with very good contiguity and partially separated phases. Our assembly has no ambiguous bases
277 with a well applicable protein-encoding gene annotation. This assembly unraveled the genomic
278 structure of *D. dumetorum* to a large extent and will serve as a reference genome sequence for yam
279 breeding by helping to identify and develop molecular markers associated with relevant agronomic
280 traits, and to understand the evolutionary history of *D. dumetorum* and yam species in general.
281

282 **Supplementary Materials**: The following are available online:
283 File S1: Sequencing overview with ENA identifiers of runs.
284 File S2: Results of BUSCO analysis of different assembly versions.
285 File S3: General statistics of different assembly versions.
286 File S4: Comparison of different gene prediction approaches.
287 File S5: Orthogroups of predicted peptides of D. rotundata and D. dumetorum.
288 File S6: Functional annotation of predicted genes in the D. dumetorum genome sequence.
289 File S7: Average short read mapping coverage of predicted genes in the D. dumetorum genome sequence.
290

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

304

## References

306 1.     Viruel, J.; Forest, F.; Paun, O.; Chase, M.W.; Devey, D.; Couto, R.S.; Segarra-Moragues, J.G.;
307     Catalan, P.; Wilkin, P. A nuclear Xdh phylogenetic analysis of yams (Dioscorea
308     Dioscoreaceae) congruent with plastid trees reveals a new Neotropical lineage. *Botanical*
309     *Journal of the Linnean Society* **2018**, 1-15.
310 2.     Sefa-Dedeh, S.; E.O., A. Biochemical and textural changes in trifoliate yam Dioscorea
311     dumetorum tubers after harvest. *Food Chemistry* **2002**, *79*, 27-40.
312 3.     Alozie, Y.E.; Akpanabiatu, M.; Eyong, E.U.; Umoh, I.B.; Alozie, G. Amino Acid Composition
313     of Dioscorea dumetorum Varities. *Pakistan Journal of Nutrition* **2009**, *8*, 103-105.
314 4.     Ferede, R.; Maziya-Dixon, B.; Alamu, O.E.; Asiedu, R. Identification and quantification of
315     major carotenoids of deep yellow-fleshed yam (tropical Dioscorea dumetorum). *Journal of*
316     *Food, Agriculture & Environment* **2010**, *8*, 160-166.

317   5.    Nimenibo-Uadia, R.; Oriakhi, A. Proximate, Mineral and Phytochemical Composition of
318         Dioscorea dumetorium Pax. *J. Appl. Sci. Environ. Manage.* **2017**, *21*, 771-774.
319   6.    Iwu, M.M.; Okunji, C.O.; Ohiaeri, G.O.; Akah, P.; Corley, D.; Tempesta, M.S. Hypoglycaemic
320         activity of dioscoretine from tubers of Dioscorea dumetorum in normal and alloxan diabetic
321         rabbits. *Planta Medica* **1990**, *56*, 264-267.
322   7.    Siadjeu, C.; Panyoo, E.A.; Toukam, G.M.S; Bell, J.M.; Nono, B.; Medoua, G.N. Influence of
323         Cultivar on the Postharvest Hardening of Trifoliate Yam (Dioscorea dumetorum) Tubers.
324         *Hindawi* **2016**, *16*.
325   8.    Tamiru, M.; Natsume, S.; Takagi, H.; White, B.; Yaegashi, H.; Shimizu, M.; Yoshida, K.;
326         Uemura, A.; Oikawa, K.; Abe, A., et al. Genome sequencing of the staple food crop white
327         Guinea yam enables the development of a molecular marker for sex determination. *BMC*
328         *Biology* **2017**, *15*, 86.
329   9.    Magwe-Tindo, J.; Wieringa, J.J.; Sonke, B.; Zapfack, L.; Vigouroux, Y.; Couvreur, T.L.P.;
330         Scarcelli, N. Complete plastome sequence of 14 African yam species (Dioscorea spp.).
331         *Mitochondrial DNA Part B* **2019**.
332   10.   Siadjeu, C.; Mayland-Quellhorst, E.; Albach, D.C. Genetic diversity and population structure
333         of trifoliate yam (Dioscorea dumetorum Kunth) in Cameroon revealed by
334         genotyping-by-sequencing (GBS). *BMC Plant Biology* **2018**, *18*, 359.
335   11.   Rosso, M.G.; Li, Y.; Strizhov, N.; Reiss, B.; Dekker, K.; Weisshaar, B. An *Arabidopsis thaliana*
336         T-DNA mutagenised population (GABI-Kat) for flanking sequence tag based reverse
337         genetics. *Plant Molecular Biology* **2003**, *53*, 247-259.
338   12.   Koren, S.; Walenz, B.P.; Berlin, K.; Miller, J.R.; Bergman, N.H.; Phillippy, A.M. Canu:
339         scalable and accurate long-read assembly via adaptive k-mer weighting and repeat
340         separation. *Genome Research* **2017**, *27*, 722-736.
341   13.   Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **2018**, *34*,
342         3094-3100.
343   14.   Vaser, R.; Sović, I.; Nagarajan, N.; Šikić, M. Fast and accurate de novo genome assembly
344         from long uncorrected reads. *Genome Research* **2017**, *27*, 737-746.
345   15.   Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
346         *arXiv* **2013**, 1303.3997v1302 (Preprint posted May 1326, 2013).
347   16.   Walker, B.J.; Abeel, T.; Shea, T.; Priest, M.; Abouelliel, A.; Sakthikumar, S.; Cuomo, C.A.;
348         Zeng, Q.; Wortman, J.; Young, S.K., et al. Pilon: an integrated tool for comprehensive
349         microbial variant detection and genome assembly improvement. *PLoS ONE* **2014**, *9*, e112963.
350   17.   Pucker, B.; Holtgräwe, D.; Rosleff Sörensen, T.; Stracke, R.; Viehöver, P.; Weisshaar, B. A De
351         Novo Genome Sequence Assembly of the Arabidopsis thaliana Accession Niederzenz-1
352         Displays Presence/Absence Variation and Strong Synteny. *PLoS ONE* **2016**, *11*, e0164321.
353   18.   Pucker, B.; Holtgräwe, D.; Weisshaar, B. Consideration of non-canonical splice sites
354         improves gene prediction on the Arabidopsis thaliana Niederzenz-1 genome sequence. *BMC*
355         *Research Notes* **2017**, *10*, 667.
356   19.   Simão, F.A.; Waterhouse, R.M.; Ioannidis, P.; Kriventseva, E.V.; Zdobnov, E.M. BUSCO:
357         assessing genome assembly and annotation completeness with single-copy orthologs.
358         *Bioinformatics* **2015**, *31*, 3210-3212.

359  20.  Keller, O.; Kollmar, M.; Stanke, M.; Waack, S. A novel hybrid gene prediction method
360       employing protein multiple sequence alignments. *Bioinformatics* **2011**, *27*, 757-763.
361  21.  Pucker, B.; Feng, T.; Brockhington, S. Next generation sequencing to investigate genomic
362       diversity in Caryophyllales. *bioRxiv* **2019**, , doi:10.1101/646133 (Preprint posted 2019-07-27).
363  22.  Altschul, S.F.; Gish, W.; Miller, W.; Myers, E.W.; Lipman, D.J. Basic local alignment search
364       tool. *Journal of Molecular Biology* **1990**, *215*, 403-410.
365  23.  Jones, P.; Binns, D.; Chang, H.Y.; Fraser, M.; Li, W.; McAnulla, C.; McWilliam, H.; Maslen, J.;
366       Mitchell, A.; Nuka, G., et al. InterProScan 5: genome-scale protein function classification.
367       *Bioinformatics* **2014**, *30*, 1236-1240, doi:10.1093/bioinformatics/btu031.
368  24.  Cheng, C.Y.; Krishnakumar, V.; Chan, A.; Thibaud-Nissen, F.; Schobel, S.; Town, C.D.
369       Araport11: a complete reannotation of the Arabidopsis thaliana reference genome. *The Plant
370       Journal* **2017**, 789-804.
371  25.  Lowe, T.M.; Eddy, S.R. tRNAscan-SE: a program for improved detection of transfer RNA
372       genes in genomic sequence. *Nucleic Acids Research* **1997**, *25*, 955-964.
373  26.  Chan, P.P.; Lowe, T.M. tRNAscan-SE: Searching for tRNA genes in genomic sequences. In
374       *Gene Prediction: Methods and Protocols*, 2019/04/26 ed.; Kollmar, M., Ed. Springer New York:
375       New York, 2019; Vol. 1962, pp. 1-14.
376  27.  Nawrocki, E.P.; Eddy, S.R. Infernal 1.1: 100-fold faster RNA homology searches.
377       *Bioinformatics* **2013**, *29*, 2933-2935.
378  28.  Kalvari, I.; Argasinska, J.; Quinones-Olvera, N.; Nawrocki, E.P.; Rivas, E.; Eddy, S.R.;
379       Bateman, A.; Finn, R.D.; Petrov, A.I. Rfam 13.0: shifting to a genome-centric resource for
380       non-coding RNA families. *Nucleic Acids Research* **2018**, *46*, D335-D342.
381  29.  Pucker, B.; Ruckert, C.; Stracke, R.; Viehover, P.; Kalinowski, J.; Weisshaar, B. Twenty-Five
382       Years of Propagation in Suspension Cell Culture Results in Substantial Alterations of the
383       Arabidopsis Thaliana Genome. *Genes* **2019**, *10*, 671, doi:10.3390/genes10090671.
384  30.  Baasner, J.S.; Howard, D.; Pucker, B. Influence of neighboring small sequence variants on
385       functional impact prediction. *bioRxiv* **2019**, , doi:10.1101/596718 (Preprint posted 2019-06-13).
386  31.  Emms, D.M.; Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
387       genomics. *Genome Biology* **2019**, *20*, 238, doi:10.1186/s13059-019-1832-y.
388  32.  Ruggieri, V.; Alexiou, K.G.; Morata, J.; Argyris, J.; Pujol, M.; Yano, R.; Nonaka, S.; Ezura, H.;
389       Latrasse, D.; Boualem, A., et al. An improved assembly and annotation of the melon
390       (Cucumis melo L.) reference genome. *Scientic Reports* **2018**, *8*, 8088,
391       doi:10.1038/s41598-018-26416-2.
392  33.  Mignouna, H.D.; Abang, M.M.; Asiedu, R. Harnessing modern biotechnology for tropical
393       tuber crop improvement: Yam (Dioscorea spp.) molecular breeding. *African Journal of
394       Biotechnology* **2003**, *2*, 12.
395  34.  Mignouna, H.D.; Abang, M.M.; Asiedu, R. Genomics of Yams, a Common Source of Food
396       and Medicine in the Tropics. In *Genomics of Tropical Crop Plants*, Moore, P.H., Ming, R., Eds.
397       2008; 10.1007/978-0-387-71219-2_23pp. 549-570.
398  35.  Michael, T.P.; Jupe, F.; Bemm, F.; Motley, S.T.; Sandoval, J.P.; Lanz, C.; Loudet, O.; Weigel,
399       D.; Ecker, J.R. High contiguity Arabidopsis thaliana genome assembly with a single
400       nanopore flow cell. *Nature Communications* **2018**, *9*, 541.

401   36.   Paajanen, P.; Kettleborough, G.; Lopez-Girona, E.; Giolai, M.; Heavens, D.; Baker, D.; Lister,
402         A.; Cugliandolo, F.; Wilde, G.; Hein, I., et al. A critical comparison of technologies for a plant
403         genome sequencing project. *Gigascience* **2019**, *8*, doi:10.1093/gigascience/giy163.
404   37.   Salmela, L.; Walve, R.; Rivals, E.; Ukkonen, E. Accurate self-correction of errors in long reads
405         using de Bruijn graphs. *Bioinformatics* **2017**, *33*, 799-806, doi:10.1093/bioinformatics/btw321.
406   38.   Miege, J. Nombres chromosomiques et répartition géographique de quelques plantes
407         tropicales et équatoriales. *Revue de Cytologie et de Biologie Végétales* **1954**, *15*, 312-348.
408   39.   Knuth, R. Dioscoreaceae. In *Das Pflanzenreich*, Engelr, A., Ed. Engelmann, W.: Leipzig, 1924.
409   40.   Wendel, J.F.; Jackson, S.A.; Meyers, B.C.; Wing, R.A. Evolution of plant genome architecture.
410         *Genome Biology* **2016**, *17*, 37, doi:10.1186/s13059-016-0908-1.
411   41.   Pucker, B.; Brockington, S.F. Genome-wide analyses supported by RNA-Seq reveal
412         non-canonical splice sites in plant genomes. *BMC Genomics* **2018**, *19*, 980.
413   42.   Pucker, B.; Holtgräwe, D.; Stadermann, K.B.; Frey, K.; Huettel, B.; Reinhardt, R.; Weisshaar,
414         B. A chromosome-level sequence assembly reveals the structure of the Arabidopsis thaliana
415         Nd-1 genome and its gene set. *PLoS One* **2019**, *14*, e0216233.
416   43.   Barton, H. Yams: Origins and Development. **2014**.
417   44.   Treche, S.; Delpeuch, F. Physiologie Vegetale - Mise en evidence de l'apparition d'un
418         epaississement membranaire dans le parenchyme des tubercules de Dioscorea dumetorum
419         au cours de la conservation. *C. R. Acad. Sc. Paris* **1979**, *288*.
420   45.   Price, E.J.; Wilkin, P.; Sarasan, V.; Fraser, P.D. Metabolite profiling of Dioscorea (yam)
421         species reveals underutilised biodiversity and renewable sources for high-value
422         compounds. *Scientic Reports* **2016**, *6*.

423