

1 **Genome sequencing of *Musa acuminata* Dwarf Cavendish reveals a chromosome arm**
2 **duplication**

3

4

5 Mareike Busche⁺, Boas Pucker⁺, Prisca Viehöver, Bernd Weisshaar, Ralf Stracke*

6

7 Genetics and Genomics of Plants, Center for Biotechnology (CeBiTec), Bielefeld University,
8 Sequenz 1, 33615 Bielefeld, Germany

9

10 MB: mbusche@cebitec.uni-bielefeld.de, 0000-0002-2114-7613

11 BP: bpucker@cebitec.uni-bielefeld.de, 0000-0002-3321-7471

12 PV: viehoeve@cebitec.uni-bielefeld.de, 0000-0003-3286-4121

13 BW: weisshaa@cebitec.uni-bielefeld.de, 0000-0002-7635-3473

14 RS: ralf.stracke@uni-bielefeld.de, 0000-0002-9261-2279

15

16 + these authors contributed equally

17 * corresponding author

18

19

20 **Abstract**

21 Different *Musa* species, subspecies, and cultivars are currently investigated to reveal their
22 genomic diversity. Here, we compare the *Musa acuminata* cultivar Dwarf Cavendish against the
23 previously released Pahang assembly. Numerous small sequence variants were detected and
24 the ploidy of the cultivar presented here was determined as triploid based on sequence variant
25 frequencies. Illumina sequencing also revealed a duplication of a large segment of chromosome
26 2 in the genome of the cultivar studied. Comparison against previously sequenced cultivars
27 provided evidence that this duplication is unique to Dwarf Cavendish. Although no functional
28 relevance of this duplication was identified, this example shows the potential of plants to tolerate
29 such aneuploidies.

30

31 **Introduction**

32 Bananas (*Musa*) are monocotyledonous perennial plants. The edible fruit (botanically a berry)
33 belongs to the most popular fruits in the world. In 2016, about 5.5 million hectares of land were
34 used for the production of more than 112 million tons of bananas [1]. The majority of bananas
35 were grown in Africa, Latin America, and Asia where they offer employment opportunities and
36 are important export commodities [1]. Furthermore, with an annual *per capita* consumption of
37 more than 200 kg in Rwanda and more than 100 kg in Angola, bananas provide food security in
38 developing countries [1,2]. While plantains or cooking bananas are commonly eaten as a staple
39 food in Africa and Latin America, the softer and sweeter dessert bananas are popular in Europe
40 and Northern America. Between 1998 and 2000, around 47 % of the world banana production
41 and the majority of the dessert banana production relied on the Cavendish subgroup of cultivars
42 [2].

43 Even though only Cavendish bananas are traded internationally, numerous cultivars are used
44 for local consumption in Africa and South-east Asia. Bananas went through a long
45 domestication process which started at least 7,000 years ago [3]. The first step towards edible
46 bananas was interspecific hybridisation between subspecies from different regions, which
47 caused incorrect meiosis and diploid gametes [4]. The diversity of edible triploids resulted from
48 human selection and triploidization of *Musa acuminata* and *Musa balbisiana* cultivars [4].

49 These exciting insights into the evolution of bananas were revealed by the analysis of genome
50 sequences. Technological advances boosted sequencing capacities and allowed the
51 (re-)sequencing of genomes from multiple subspecies and cultivars. *M. acuminata* can be
52 divided into several subspecies and cultivars. The first *M. acuminata* (DH Pahang) genome
53 sequence has been published in 2012 [5], many more genomes have been sequenced recently
54 including: banksii, burmannica, zebrina [6], malaccensis (SRR8989632, SRR6996493), Baxijiao
55 (SRR6996491, SRR6996491), Sucrier : Pisang_Mas (SRR6996492). Additionally, the genome
56 sequences of other *Musa* species, *M. balbisiana* [7], *M. itinerans* [8], and *M. schizocarpa* [9],
57 have already been published.

58

59 Here we report about our investigation of the genome of *M. acuminata* Dwarf Cavendish which
60 revealed multiple large scale deletions compared to the DH Pahang v2 reference assembly.
61 Moreover, we identified an increased copy number of the southern arm of chromosome 2,
62 indicating that this region was duplicated in one haplophase.

63

64 **Materials and methods**

65 **Plant material and DNA extraction.** *Musa acuminata* Dwarf Cavendish tissue culture
66 seedlings were obtained from FUTURE EXOTICS/SolarTek (Düsseldorf, Germany) (Figure 1).
67 Plants were grown under natural daylight at 21°C. Genomic DNA was isolated from leaves
68 following the protocol of Dellaporta *et al.* [10].



69

70 **Figure 1: *M. acuminata* Dwarf Cavendish plant.**

71

72 **Library preparation and sequencing.** Genomic DNA was subjected to sequencing library
73 preparation via the TrueSeq v2 protocol as previously described [11]. Paired-end sequencing
74 was performed on an Illumina HiSeq1500 and NextSeq500, respectively, resulting in 2x250 nt
75 and 2x154 nt read data sets.

76

77 **Read mapping, variant calling, and variant annotation.** All reads were mapped to the DH
78 Pahang v2 reference genome sequence via BWA-MEM v0.7 [12]. GATK v3.8 [13,14] was
79 deployed to identify small sequence variants based on this read mapping. The resulting variant
80 set was subjected to SnpEff [15] and NAVIP [16] to assign predictions about the functional
81 impact to all variants. Variants with disruptive effects were selected using a customized Python
82 script as described earlier [11].

83 The genome-wide distribution of small sequence variants was assessed for small nucleotide
84 variants (SNVs) and InDels based on previously developed scripts [16]. The length distribution
85 of InDels inside coding sequences was compared to the length distribution of InDels outside
86 coding sequences using a customized Python script [11].

87

88 **De novo genome assembly.** Trimmomatic v0.38 [17] was applied as previously described [11]
89 to remove low quality sequences and remaining adapter sequences. Different sets of trimmed
90 reads were subjected to SOAPdenovo2 for assembly using optimized parameters [18].
91 Resulting assemblies were evaluated using previously described criteria [18] including BUSCO
92 v3 [19] and polished by removing potential contaminations and adapters as described before
93 [18]. The DH Pahang v2 assembly [5,20] was used in the contamination detection process to
94 distinguish between *bona fide* banana contigs and sequences of unknown origin. Contigs with
95 high sequence similarity to non-plant sequences were removed as previously described [18].
96 Surviving contigs were sorted based on the DH Pahang v2 reference genome sequence and
97 joint into pseudochromosomes to facilitate downstream analyses.

98

99 **Data availability.** Sequencing datasets were submitted to EBI (ERR3412983, ERR3412984,
100 ERR3413471, ERR3413472, ERR3413473, ERR3413474). Python scripts are freely available
101 on github (<https://github.com/bpucker/banana>).

102

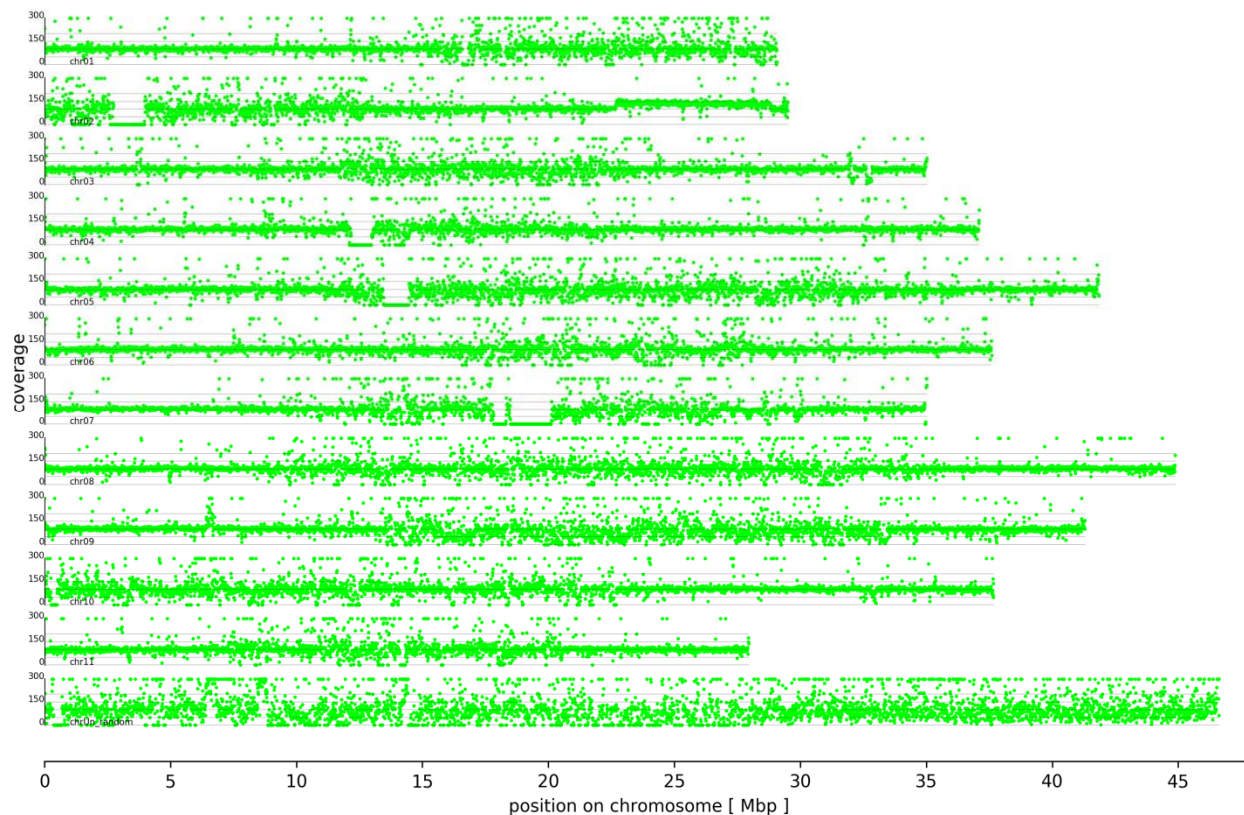
103

104 **Results and discussion**

105 **Structural variants.** The probably most remarkable difference between the Dwarf Cavendish
106 and Pahang genome sequence is the amplification of the southern region of chromosome 2
107 (Figure 2). Such a duplication was not observed in any of the other available sequencing data
108 sets (File S1, File S2). There are multiple deletions and insertions which are not randomly
109 distributed over the genome. An increased average coverage in the south of chromosome 2
110 indicates that a substantial fraction of this chromosome has been duplicated. Apparently, read
111 mapping indicates at least four large scale deletions in Dwarf Cavendish compared to Pahang
112 v2 on chromosomes 2, 4, 5 and 7 (Figure 2). However, analysis of the underlying sequence
113 revealed long stretches of ambiguous bases as the cause for these low coverage regions.

114

115



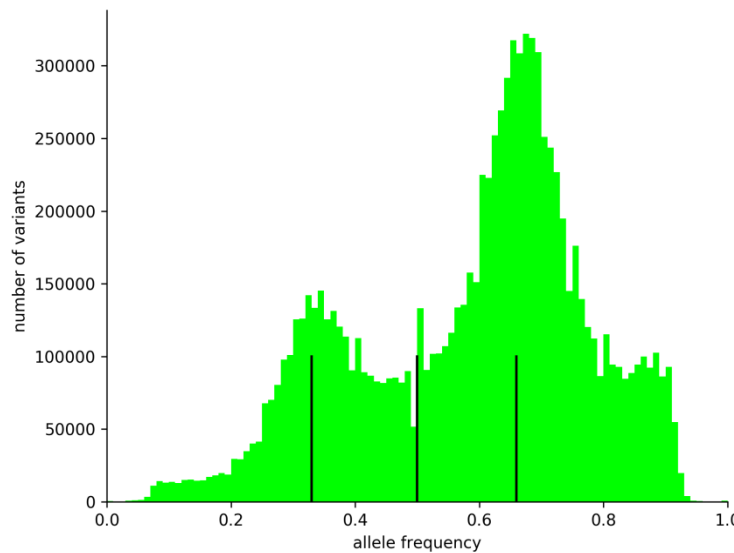
116

117 **Figure 2: Coverage distribution.** Chromosomes are ordered by increasing number with the
118 north end on the left hand side. Mapping of *M. acuminata* Dwarf Cavendish reads against the
119 DH Pahang v2 reference sequence revealed a tetraploid region in the southern part of
120 chromosome 2 in Dwarf Cavendish. Apparent large scale deletions are technical artifacts
121 caused by large stretches of ambiguous bases in the Pahang assembly that cannot be covered
122 by reads.

123

124 **Ploidy of *M. acuminata* Dwarf Cavendish.** Based on the coverage of small sequence variants,
125 the ploidy of Dwarf Cavendish was identified as triploid (Figure 3). Many heterozygous variants
126 display a frequency of the reference allele close to 0.33 or close to 0.66. This can be explained
127 by two copies of the reference allele and one copy of a different allele or *vice versa*. Deviation
128 from the precise values can be explained by random fluctuation. Since the peak around 0.66 is
129 substantially higher than the peak around 0.33, it is possible that two haplophases are very
130 similar to the reference, while one haplophase is different. This hypothesis could be tested in a
131 high continuity phased assembly.

132



133

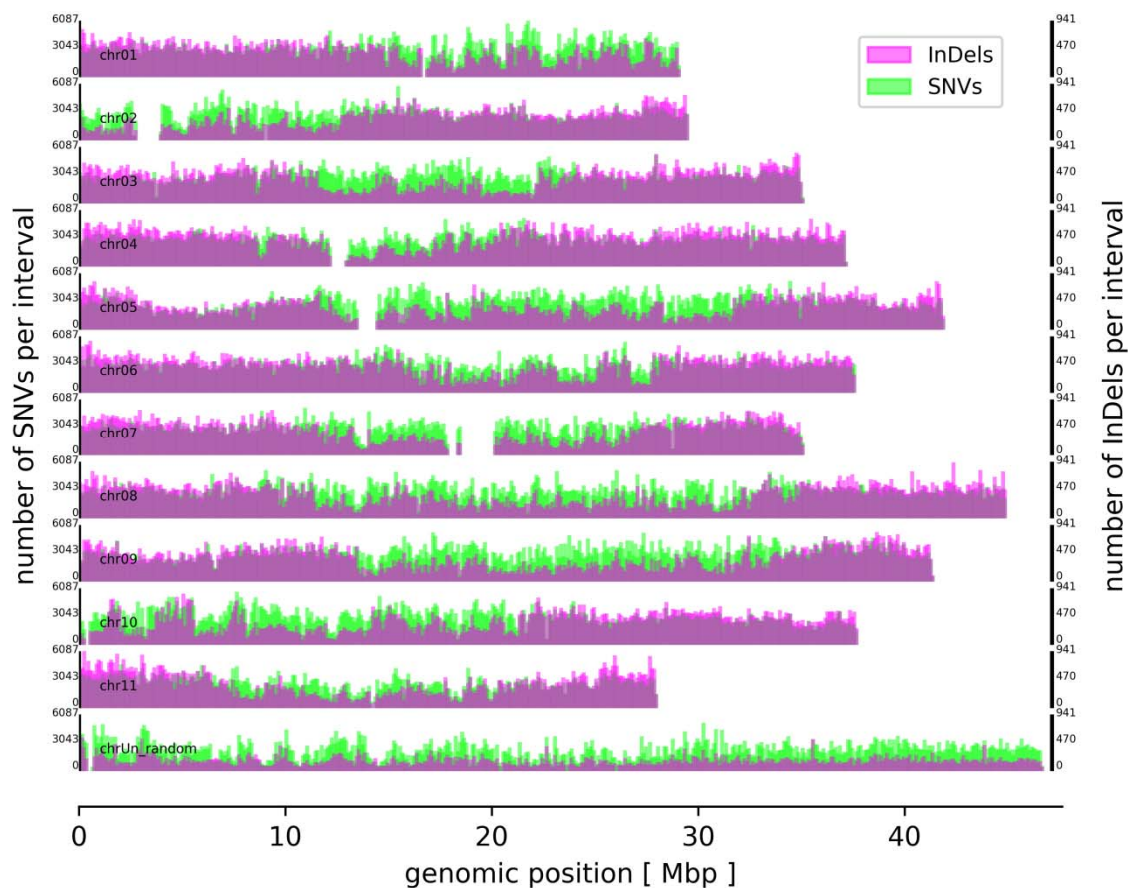
134 **Figure 3: Allele frequency histogram.** A mapping of Illumina reads against the Pahang v2
135 reference sequence was used for the identification of small sequence variants. The frequencies
136 of the reference allele at small variant positions are displayed here, excluding those positions at
137 which the Pahang v2 reference sequence deviates from an invariant sequence position of Dwarf
138 Cavendish. Black vertical lines indicate 0.33, 0.5, and 0.66, respectively.

139

140 **Genome-wide distribution of small sequence variants.** In total, 12,012,019 single nucleotide
141 variants and 1,546,394 InDels of up to 894 bp were identified between the Dwarf Cavendish
142 reads and the Pahang v2 assembly ([https://docs.cebitec.uni-
143 bielefeld.de/s/grCccg2p3WPScWZ](https://docs.cebitec.uni-bielefeld.de/s/grCccg2p3WPScWZ)). As previously observed in other re-sequencing studies [11],
144 the number of SNVs exceeds the number of InDels substantially. Moreover, InDels are more
145 frequent outside of annotated coding regions. Inside coding regions, InDels show an increased
146 proportion of lengths which are divisible by 3, because no deleterious frameshifts are caused.

147 SnpEff predicted 4,384 premature stop codons, 3,345 lost stop codons, and 8,296 frameshifts
148 based on this variant set (File S3). Even given the larger genome size, these numbers are
149 substantially higher than high impact variant numbers observed in re-sequencing studies of
150 homozygous species before [11,21]. One explanation could be the presence of three alleles for
151 each locus leading to compensation of disrupted alleles. Since banana plants are propagated
152 vegetatively, breeders do not suffer inbreeding depressions.

153



154

155 **Figure 4: Genome-wide distribution of small sequence variants.** Single nucleotide variants
 156 (SNVs, green) and small InDels (magenta) distinguish Dwarf Cavendish from Pahang. Variants
 157 were counted in 100 kb windows and are displayed on two different y-axes to allow maximal
 158 resolution [11].

159

160 **De novo genome assembly.** To facilitate wet lab applications like oligonucleotide design and
 161 validation of amplicons, the genome sequence of Dwarf Cavendish was assembled *de novo*.
 162 The assembly comprises 404,889 contigs with an N50 of 4.7 kb (Table 1). Differences between
 163 the three haplophases are one possible explanation for the low assembly contiguity. Since the
 164 short reads are insufficient to resolve entire haplophases, the assembly size is only slightly
 165 exceeding the size of one haplotype. Due to the low contiguity of this assembly and only
 166 minimal above 50% complete BUSCOs (Benchmarking Universal Single-Copy Orthologs) [19],
 167 annotation was omitted.

168

169 **Table 1: *M. acuminata* Dwarf Cavendish *de novo* genome assembly statistics.**

Parameter	Value
-----------	-------

Number of contigs	404,889
Maximal contig length	240,314 bp
Assembly size	1,068,694,790 bp
GC content	38.87 %
N50	4,709
N90	1,006

170

171

172 **Acknowledgments**

173 We thank Joachim Weber for great technical assistance.

174

175 **Author contributions**

176 BP, MB and RS planned the experiment. PV did the library preparation and sequencing. BP
177 performed bioinformatic analyses. MB and BP wrote the initial draft. MB, BP, BW and RS
178 revised the manuscript. All authors read and approved the final manuscript version.

179

180

181 **References**

- 182 1. FAO. FAOSTAT [Internet]. [cited 23 Jun 2019]. Available:
183 <http://www.fao.org/faostat/en/#data/QC>
- 184 2. Arias P, Dankers C, Liu P, Pilkauskas P. The World Banana Economy, 1985-2002. Food
185 and Agriculture Organisation (FAO) of the United Nations 2003. Available:
186 <http://www.fao.org/3/y5102e/y5102e00.htm>
- 187 3. Denham TP, Haberle SG, Lentfer C, Fullagar R, Field J, Therin M, et al. Origins of
188 Agriculture at Kuk Swamp in the Highlands of New Guinea. *Science*. 2003;301: 189–193.
189 doi:10.1126/science.1085255
- 190 4. Perrier X, Langhe ED, Donohue M, Lentfer C, Vrydaghs L, Bakry F, et al. Multidisciplinary
191 perspectives on banana (*Musa* spp.) domestication. *Proc Natl Acad Sci*. 2011;108: 11311–
192 11318. doi:10.1073/pnas.1102001108
- 193 5. D'Hont A, Denoeud F, Aury J-M, Baurens F-C, Carreel F, Garsmeur O, et al. The banana
194 (*Musa acuminata*) genome and the evolution of monocotyledonous plants. *Nature*.
195 2012;488: 213–217. doi:10.1038/nature11241
- 196 6. Rouard M, Droc G, Martin G, Sardos J, Hueber Y, Guignon V, et al. Three New Genome
197 Assemblies Support a Rapid Radiation in *Musa acuminata* (Wild Banana). *Genome Biol*
198 *Evol*. 2018;10: 3129–3140. doi:10.1093/gbe/evy227

- 199 7. Davey MW, Gudimella R, Harikrishna JA, Sin LW, Khalid N, Keulemans J. “A draft *Musa*
200 *balbisiana* genome sequence for molecular genetics in polyploid, inter- and intra-specific
201 *Musa* hybrids.” *BMC Genomics*. 2013;14: 683. doi:10.1186/1471-2164-14-683
- 202 8. Wu W, Yang Y-L, He W-M, Rouard M, Li W-M, Xu M, et al. Whole genome sequencing of
203 a banana wild relative *Musa itinerans* provides insights into lineage-specific diversification
204 of the *Musa* genus. *Sci Rep*. 2016;6: 31586. doi:10.1038/srep31586
- 205 9. Belser C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, et al. Chromosome-
206 scale assemblies of plant genomes using nanopore long reads and optical maps. *Nat*
207 *Plants*. 2018;4: 879. doi:10.1038/s41477-018-0289-4
- 208 10. Dellaporta SL, Wood J, Hicks JB. A plant DNA miniprep: Version II. *Plant Mol Biol*
209 *Report*. 1983;1: 19–21. doi:10.1007/BF02712670
- 210 11. Pucker B, Holtgräwe D, Sörensen TR, Stracke R, Viehöver P, Weisshaar B. A *De Novo*
211 *Genome Sequence Assembly of the Arabidopsis thaliana* Accession Niederzenz-1
212 Displays Presence/Absence Variation and Strong Synteny. *PLOS ONE*. 2016;11:
213 e0164321. doi:10.1371/journal.pone.0164321
- 214 12. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
215 *ArXiv13033997 Q-Bio*. 2013; Available: <http://arxiv.org/abs/1303.3997>
- 216 13. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
217 *Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA*
218 *sequencing data*. *Genome Res*. 2010;20: 1297–1303. doi:10.1101/gr.107524.110
- 219 14. Auwera GAV der, Carneiro MO, Hartl C, Poplin R, Angel G del, Levy Moonshine A, et al.
220 *From FastQ Data to High-Confidence Variant Calls: The Genome Analysis Toolkit Best*
221 *Practices Pipeline*. *Curr Protoc Bioinforma*. 2013;43: 11.10.1-11.10.33.
222 doi:10.1002/0471250953.bi1110s43
- 223 15. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for
224 annotating and predicting the effects of single nucleotide polymorphisms, SnpEff. *Fly*
225 (Austin). 2012;6: 80–92. doi:10.4161/fly.19695
- 226 16. Baasner J-S, Howard D, Pucker B. Influence of neighboring small sequence variants on
227 functional impact prediction. *bioRxiv*. 2019; 596718. doi:10.1101/596718
- 228 17. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
229 data. *Bioinforma Oxf Engl*. 2014;30: 2114–2120. doi:10.1093/bioinformatics/btu170
- 230 18. Pucker B, Feng T, Brockington S. Next generation sequencing to investigate genomic
231 diversity in Caryophyllales. 2019; doi:<https://doi.org/10.1101/646133>
- 232 19. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
233 assessing genome assembly and annotation completeness with single-copy orthologs.
234 *Bioinforma Oxf Engl*. 2015;31: 3210–3212. doi:10.1093/bioinformatics/btv351

- 235 20. Martin G, Baurens F-C, Droc G, Rouard M, Cenci A, Kilian A, et al. Improvement of the
236 banana “*Musa acuminata*” reference sequence using NGS data and semi-automated
237 bioinformatics methods. *BMC Genomics*. 2016;17: 243. doi:10.1186/s12864-016-2579-4
- 238 21. Xu Y-C, Niu X-M, Li X-X, He W, Chen J-F, Zou Y-P, et al. Adaptation and Phenotypic
239 Diversification in *Arabidopsis* through Loss-of-Function Mutations in Protein-Coding
240 Genes. *Plant Cell*. 2019;31: 1012–1025. doi:10.1105/tpc.18.00791

241

242

243 **Additional files**

244 **File S1.** List of public sequence read samples used for comparison against Dwarf Cavendish
245 based on the DH Pahang reference.

246 **File S2.** Coverage plots of public sequence read samples for comparison against Dwarf
247 Cavendish based on the DH Pahang reference. Order of samples is: *Musa acuminata*,
248 AYP_BOSN_r1, *Musa banksia*, *Musa burmannica*, Cavendish BaXiJiao, Gros Michel, *Musa*
249 *malaccensis*, Sucrier : Pisang Mas, sucrier pisang mas 1998-2307, *Musa zebrina*,
250 Pisang_Klutuk_Wulung, *Musa itinerans*, *Musa schizocarpa*

251 **File S3.** Selected high impact variants between Dwarf Cavendish and DH Pahang with effects
252 predicted by SnpEff.