

# 1 25 years of propagation in suspension cell culture results 2 in substantial alterations of the *Arabidopsis thaliana* 3 genome

4 Boas Pucker<sup>1,\*</sup>, Christian Rückert<sup>2</sup>, Ralf Stracke<sup>1</sup>, Prisca Viehöver<sup>1</sup>, Jörn Kalinowski<sup>2</sup> and Bernd  
5 Weisshaar<sup>1</sup>

6 1 Genetics and Genomics of Plants, Faculty of Biology, Center for Biotechnology (CeBiTec), Bielefeld University,  
7 Sequenz 1, 33615 Bielefeld, NRW, Germany

8 2 Microbial Genomics and Biotechnology, Center for Biotechnology (CeBiTec), Bielefeld University, Sequenz 1, 33615  
9 Bielefeld, NRW, Germany

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

11 CR: cruecker@cebitec.uni-bielefeld.de, 0000-0002-9722-4435

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

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

14 JK: joern@cebitec.uni-bielefeld.de, 0000-0002-9052-1998

15 BW: weissshaa@cebitec.uni-bielefeld.de, 0000-0002-7635-3473

16 \* correspondence: bpucker@cebitec.uni-bielefeld.de

17 **Abstract:** *Arabidopsis thaliana* is one of the best studied plant model organisms. Besides cultivation in  
18 greenhouses, cells of this plant can also be propagated in suspension cell culture. At7 is one such cell line that  
19 has been established about 25 years ago. Here we report the sequencing and the analysis of the At7 genome.  
20 Large scale duplications and deletions compared to the Col-0 reference sequence were detected. The number  
21 of deletions exceeds the number of insertions thus indicating that a haploid genome size reduction is ongoing.  
22 Patterns of small sequence variants differ from the ones observed between *A. thaliana* accessions e.g. the  
23 number of single nucleotide variants matches the number of insertions/deletions. RNA-Seq analysis reveals  
24 that disrupted alleles are less frequent in the transcriptome than the native ones.

25 **Keywords:** copy number variations; variant calling; next generation sequencing; long read sequencing

26

---

## 27 1. Introduction

28 *Arabidopsis thaliana* is a small flowering plant which is distributed over the northern hemisphere and has  
29 become the model system of choice for research in plant biology. In 2000, the genome sequence was released as  
30 the first available plant genome sequence [1]. After generating this reference sequence from the accession  
31 Columbia-0 (Col-0), many other *A. thaliana* accessions were analyzed by sequencing to investigate, among  
32 many other topics, genomic diversity, local adaptation and the phylogenetic history of this species [2-4]. While  
33 most initial re-sequencing projects relied on short read mapping against the Col-0 reference sequence [5-7],  
34 technological progress enabled *de novo* genome assemblies [2,8] which reached a chromosome-level quality  
35 [9-11]. Low coverage nanopore sequencing was also applied to search for genomic differences e.g. active  
36 transposable elements (TEs) [12].

37 Due to the high value of *A. thaliana* for basic plant biology research, it is frequently grown in greenhouses  
38 under controlled and optimized conditions. Previous studies investigated the mutation rates within a single  
39 generation [6,13]. Mutational changes appear to be different between plants grown under controlled conditions  
40 and natural samples collected in the environment [6]. Another approach harnessed an *A. thaliana* population in  
41 the United States of America, which is assumed to originate from a single ancestor thus showing mutations  
42 accumulated over the last decades [14]. This study investigated modern and ancient specimens and estimated a  
43 rate of  $7.1 \cdot 10^{-9}$  substitutions per site per generation [14]. Since not all mutations are fixed during evolution, the  
44 mutation rate is higher than the substitution rate.

45 Even further away from natural conditions in the environment is the propagation of cells in suspension  
46 cultures. Cells from such cultures can easily be employed for transient transfection experiments [15]. Transient  
47 transfections of At7 protoplasts are a relatively straightforward method to study promoter structure and activity  
48 and to investigate the interactions between transcription factors and promoters of putative target genes [16,17].

49 Since most functions of plants are dispensable in suspension culture, it was expected that mutations in  
50 these dispensable genes accumulate over time due to genetic drift or even due to positive selection. We  
51 sequenced and analyzed the genome of At7 cells which have been propagated in suspension culture for about 25  
52 years and identified massive genomic changes.

## 53 **2. Methods**

### 54 *2.1. Plant material*

55 The *Arabidopsis thaliana* suspension cell culture At7 [18] is derived from hypocotyl of the reference  
56 accession Columbia (Col). This cell culture is cultivated at 26°C and 105 rpm in darkness with sugar supply on  
57 a rotary shaker. A subset of cells is transferred into fresh B5 medium at a weekly basis. For details about the  
58 propagation of this cell culture see [19]. All sequencing data sets generated in this study were submitted to the  
59 European Nucleotide Archive (ENA) as part of the study PRJEB33589 and the sample ERS3588070.

### 60 *2.2. RNA extraction and RNA-Seq*

61 Total RNA was extracted based on a previously described protocol [16]. Based on 1 µg of total RNA,  
62 sequencing libraries were constructed following the TrueSeq v2 protocol. Three biological replicates of the At7  
63 cell culture (splitted the preceding week) were processed and combined for sequencing after one week of  
64 incubation. Single end sequencing of 83 nt was performed on an Illumina NextSeq500 at the Sequencing Core  
65 Facility of the Center for Biotechnology (CeBiTec) at Bielefeld University (ERR3444576, ERR3444577).

### 66 *2.3. DNA extraction*

67 Cells were separated from media through filtration on a Büchner funnel. Next, cell walls were destroyed  
68 by treatment in the Precellys Evolution (QIAGEN) for 3x30 seconds at 8500 rpm with 1 mm zirconium beads  
69 (Roth). The following DNA extraction was based on a previously described cetyltrimethylammoniumbromid  
70 (CTAB) protocol [20]. To increase the purity of the DNA an additional ethanol precipitation step was added.  
71 DNA from the same extraction was used for sequencing with different second and third generation  
72 technologies.

### 73 *2.4. Oxford Nanopore sequencing*

74 DNA was quantified via a Qubit fluorometer (Invitrogen) following suppliers' recommendation.  
75 Sequencing was performed on a total of four flow cells. The Ligation Sequencing Kit SQK-LSK109 (Oxford  
76 Nanopore Technologies, ONT) and the Rapid Sequencing Kit SQK-RAD004 (ONT) were used to prepare two  
77 DNA libraries each for sequencing based on the suppliers' protocol. Sequencing on a MinION (ONT)  
78 (ERR3445571) and GridION (ONT) (ERR3445572), respectively, was performed using R9.4 flow cells (ONT)  
79 based on the supplier's instructions. Real time basecalling for the MinION experiments was performed on a  
80 MinIT (ONT) via MinKNOW software (ONT).

### 81 *2.5. Illumina DNA sequencing*

82 Sequencing library preparation was performed as described previously [8] and sequencing took place at  
83 the Sequencing Core Facility of the CeBiTec at Bielefeld University. Paired-end sequencing of the library was  
84 done on an Illumina HiSeq1500 system in Rapid Mode to produce 2x250 nt reads (ERR3445426) and on an  
85 Illumina NextSeq500 benchtop sequencer resulting in 2x155 nt reads (ERR3445427).

### 86 *2.6. Coverage analysis*

87 Dedicated Python scripts were applied for the genome-wide coverage analysis, the classification of genes  
88 based on coverage values, and the investigation of coverage values at variant positions. These scripts are  
89 available at github: <https://github.com/bpucker/At7>. Completely deleted regions were identified through the  
90 identification of zero coverage regions in At7, i.e. regions where no At7 reads are mapped to the Col-0 reference  
91 sequence, as described before [8].

### 92 *2.7. Genome assembly*

93       ONT reads longer than 3 kb were subjected to genome assembly via Canu v1.8 assembler [21] using the  
94 following parameters: 'genomeSize=145m' 'useGrid=1' 'saveReads=true'  
95 'corMhapFilterThreshold=0.0000000002' 'ovlMerThreshold=500' 'corMhapOptions=--threshold 0.80  
96 --num-hashes 512 --num-min-matches 3 --ordered-sketch-size 1000 --ordered-kmer-size 14 --min-olap-length  
97 2000 --repeat-idf-scale 50'. In addition to Canu, Miniasm v0.3-r179 [22] and Flye v2.3.1 [23] were run to  
98 identify the best assembler for this data set as previously described [11]. The quality of assembled contigs was  
99 improved through nine successive rounds of polishing with Nanopolish v0.11 [24] harnessing the information  
100 of all ONT reads. Minimap2 v2.10-r761 [22] was used for the read mapping with the arguments  
101 '--secondary=no -ax map-ont'. Next, Illumina reads were mapped via BWA MEM v0.7.13 [25] to the polished  
102 assembly. Five successive polishing rounds with Pilon v1.22 [26] were performed based on Illumina read  
103 mappings. Previously developed Python scripts [8] were applied to check the assembly for contaminations, to  
104 remove very short contigs, and to calculate general assembly statistics.

## 105 2.8. Variant calling and annotation

106       Illumina sequencing reads of At7 were aligned to the TAIR9 reference sequence of Col-0 via BWA MEM  
107 v0.7.13 [25]. Next, GATK v.3.8 [27,28] was applied for the identification of small sequence variants as  
108 previously described [29]. In contrast to previous studies, variant positions with multiple different alleles were  
109 kept as they are biologically possible. GATK is able to identify low frequency variants and thus provides a  
110 powerful solution given the variation in ploidy between hemizygous and pentaploid. Variants were considered  
111 as heterozygous if at least 5 reads and at least 5% of all reads support an additional allele besides the dominant  
112 one. Additionally, variants were excluded if their coverage deviates from the average coverage of a 2 kb  
113 flanking region by more than 20%. SnpEff [30] was deployed to assign predictions of the functional impact to  
114 all small sequence variants. Previously developed Python scripts [8,29] were customized to investigate the  
115 distribution of variants and to check for patterns. SVIM [31] was deployed to identify large variants based on a  
116 Minimap2 mapping of all ONT reads against the Col-0 reference sequence. RNA-Seq reads were mapped via  
117 STAR [32] to the Col-0 reference genome sequence using previously described parameters [33]. Samtools [34]  
118 and bedtools [35] were used to call variants based on an RNA-Seq read mapping to the Col-0 reference  
119 sequence.

120

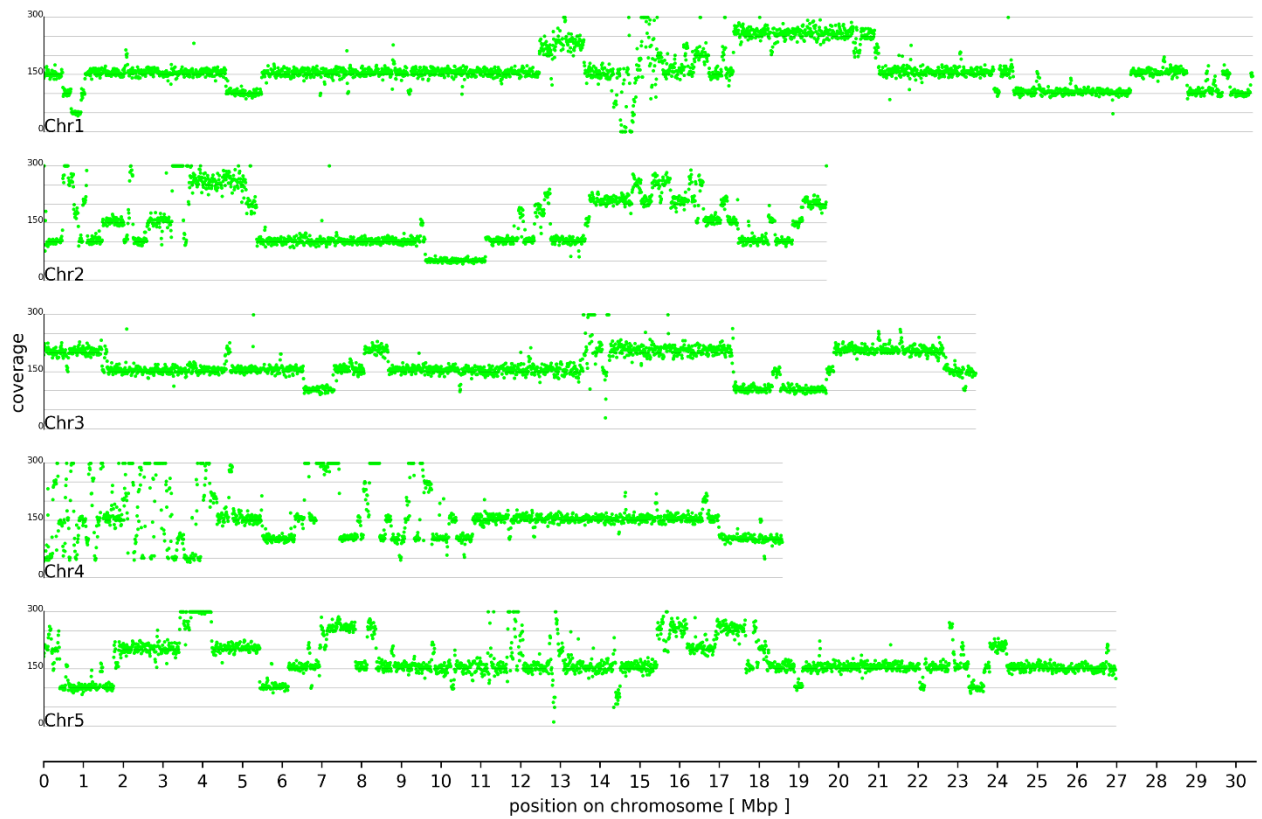
## 121 3. Results and Discussion

### 122 3.1. Ploidy differences between At7 and Col-0

123       Mapping of the At7 Illumina short reads and ONT long reads to the Col-0 reference sequence was used to  
124 assess genomic changes in At7. The coverage analysis revealed duplications and deletions of large  
125 chromosomal segments (Figure 1, File S1). Similar regional variations in ploidy between neighbouring  
126 chromosomal segments are common in immortalized insect and mammalian cell lines and tumors, where they  
127 may be an advantage to cells [36,37]. In the At7 suspension cell culture, about 5 Mbp at the northern end of  
128 chromosome 2 (Chr2) and chromosome 4 (Chr4) appear highly fragmented. In addition, regions around the  
129 centromeres are apparently fragmented, but this could be an artifact of higher repeat content and a substantial  
130 proportion of collapsed peri-centromeric and centromeric sequences. These differences in chromosomal  
131 stability seem to be consistent between plants and animals since similar observations were also reported for e.g.  
132 Chinese hamster ovary (CHO) and *Drosophila* cell lines [36,37].

133       An average coverage of 50 fold in a small region at 1 Mbp on Chr1 and in a region between 9.5 Mbp and  
134 11 Mbp on Chr2 indicates hemizyosity, i.e. there is only one copy left and the regions are hemizygous. Small  
135 blocks in the northern end of Chr4 show a similar coverage. An average coverage of 150 fold indicates that most  
136 other regions are triploid. Although the starting material from Col-0 was diploid 25 years ago, these regions  
137 represent a minority in At7. A region on Chr1 (17.5 to 21 Mbp) is present in five copies (pentaploid). There are  
138 only 116 regions larger than 100 bp which appeared of being completely deleted in At7. These observations in  
139 At7 differ from findings in CHO cells, where large chromosomal segments were found in a hemizygous state  
140 [37]. Possible explanations for the observed polyploidy regions are defective mitosis and chromosome  
141 endoreplication as previously reviewed by [38] and references therein. Similar chromosomal restructuring has

142 been reported during genome elimination in an *A. thaliana* plant with different centromer-specific H3 histones  
143 [39].



144

145 **Figure 1. Genome-wide coverage of the Col-0 reference sequence.** Hemizygous regions in At7 were revealed  
146 by a read coverage of approximately 50 fold when combining Illumina and ONT sequencing reads. Different  
147 multiples of this values can be observed revealing the presence of large scale duplications.

148

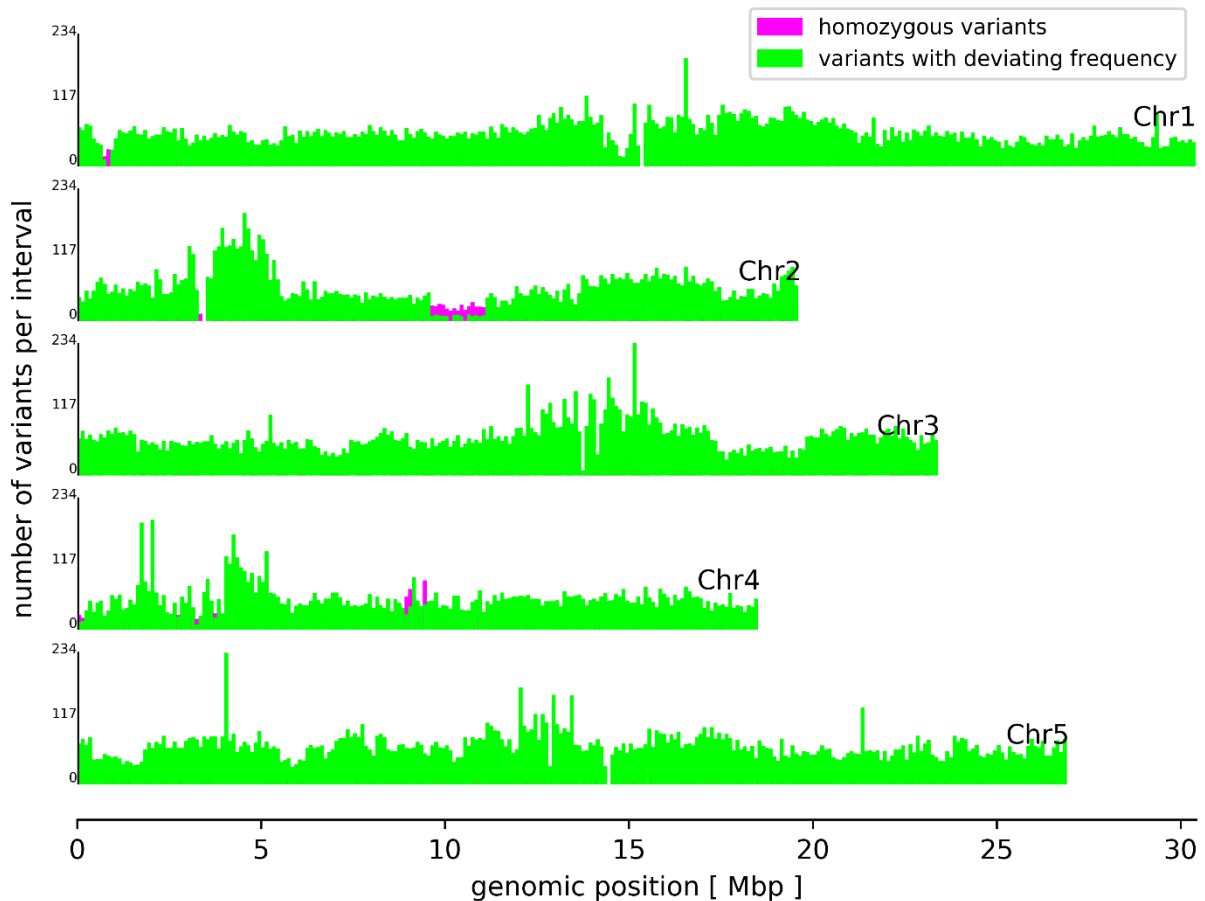
149 Besides Illumina sequencing, Oxford Nanopore Technologies' MinION and GridION sequencer were  
150 deployed to generate long reads for the detection of large structural variants. Alignments against the Col-0 reference  
151 sequence revealed 160,348 deletions and 5,902 insertions (Figure S1). Previous comparisons of natural *A. thaliana*  
152 accessions revealed equal numbers of insertions and deletions [2,8,29]. For comparison, long reads of the Nd-1  
153 accession were aligned to the Col-0 reference sequence resulting in the identification of almost equal numbers of  
154 insertions and deletions. Although different long read sequencing technologies are involved, this substantial  
155 difference between the At7 comparison to Col-0 and the Nd-1 comparison to Col-0 is unlikely to be a technical  
156 artifact. Analysis of the remaining coverage in deletions supports the validity, because most of these regions seem  
157 to be maintained in at least two haplophases (Figure S2). The excess of deletions observed here indicates that a single  
158 haplophase of the At7 genome is probably smaller than a haploid Col-0 genome. This could be due to an ongoing  
159 genome size reduction in the At7 cell culture. A subset of the At7 cells in the suspension culture is transferred to a  
160 new flask with fresh media on a weekly basis. Fast dividing cells have a higher chance of passing at least one  
161 daughter cell to the new flask. The result is a strong selection pressure for fast cell division. Although the replication  
162 of additional sequences is probably not costly in terms of evolution, the transcription and translation of dispensable  
163 genes was shown to be under negative selection in eukaryotes [40]. Artificial selection through humans is one of the  
164 strongest forces that can drive evolutionary changes [5,6,14]. The comparison of the genome sequences of  
165 *A. thaliana* and the close relative *A. lyrata* revealed randomly distributed small deletions as one major source for the  
166 genome size differences [41]. High numbers of large deletions between At7 and Col-0 indicate that the haploid  
167 genome size reduction is occurring at substantially higher speed and using different mechanisms.

168 Read mappings revealed that both, plastome and chondrome, are still detectable in At7 cells. Due to differences  
169 in the cultivation conditions and DNA extraction methods, a reliable comparison of plastome coverage to nucleom  
170 coverage between At7 and plants is not feasible. Although the loss of chloroplasts may seem likely in plant cells  
171 incubated in the dark and with sugar supply through the cultivation media, several central biosynthetic pathways like  
172 *de novo* synthesis of fatty acids [42,43] and isoprenoid biosynthesis [44] are located in the plastids [45]. This is  
173 probably the explanation why this organell is maintained in At7 cells. The coverage ratio between chondrome and  
174 plastome is about 10 times higher than observed in native plants [11,46] (Figure S3). The increased number of  
175 mitochondria could be due to the specific conditions cultured At7 cells are exposed to. Previous reports described  
176 differences in organell numbers due to development and environmental stresses [45].

177

### 178 3.2. Sequence variants and copy number variations

179 Mapping of all At7 Illumina sequencing reads against the Col-0 reference sequence and following variant  
180 calling revealed an almost equal amount of single nucleotide variants (SNVs) and small insertions/deletions  
181 (InDels) (File S2). This is in contrast to previous re-sequencing studies where natural populations displayed a  
182 substantial excess of SNVs compared to InDels [2,5,8]. In total, 127,686 small variants were identified leading  
183 to a frequency of about one variant per kbp. This frequency is lower than previously reported values of 1 variant  
184 in 200-400 bp for the comparison of natural *A. thaliana* accessions [5,8]. The sequence variants show an  
185 increased change from CG to AT, when compared to the Nd-1 vs. Col-0 comparison [8,29]. Variants in At7 are  
186 the result of spontaneous mutations which accumulate in the asexually reproducing At7 cells according to  
187 Muller's Ratchet [47] as previously reported for apomicts [48,49]. The higher ratio of InDels could indicate that  
188 more variants are deleterious than usually found between sexually reproducing accessions. This would be in  
189 agreement with the expectation based on Muller's Ratchet. A generally higher frequency of variants between  
190 accessions could be due to rare outcrossing events, which can introduce numerous variants in a single event  
191 [50]. It is even possible that several of these variants between accessions compensate each other's effect [2,29],  
192 while isolated variants in At7 are deleterious. In agreement with previous findings [5,8], the genome-wide  
193 distribution of variants between At7 and the Col-0 reference sequence showed increased variant frequencies  
194 around the centromeres (Figure 2). As previously suggested [8], this enrichment can be partly explained by a  
195 higher proportion of collapsed sequences in these regions due to the higher abundance of repeats.



196

197

198

199

**Figure 2. Genome-wide distribution of small sequence variants between At7 and Col-0.** Homozygous variants (magenta) and variants with deviating frequency were counted in genomic blocks of 100 kb on all five chromosomes of the Col-0 reference sequence.

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

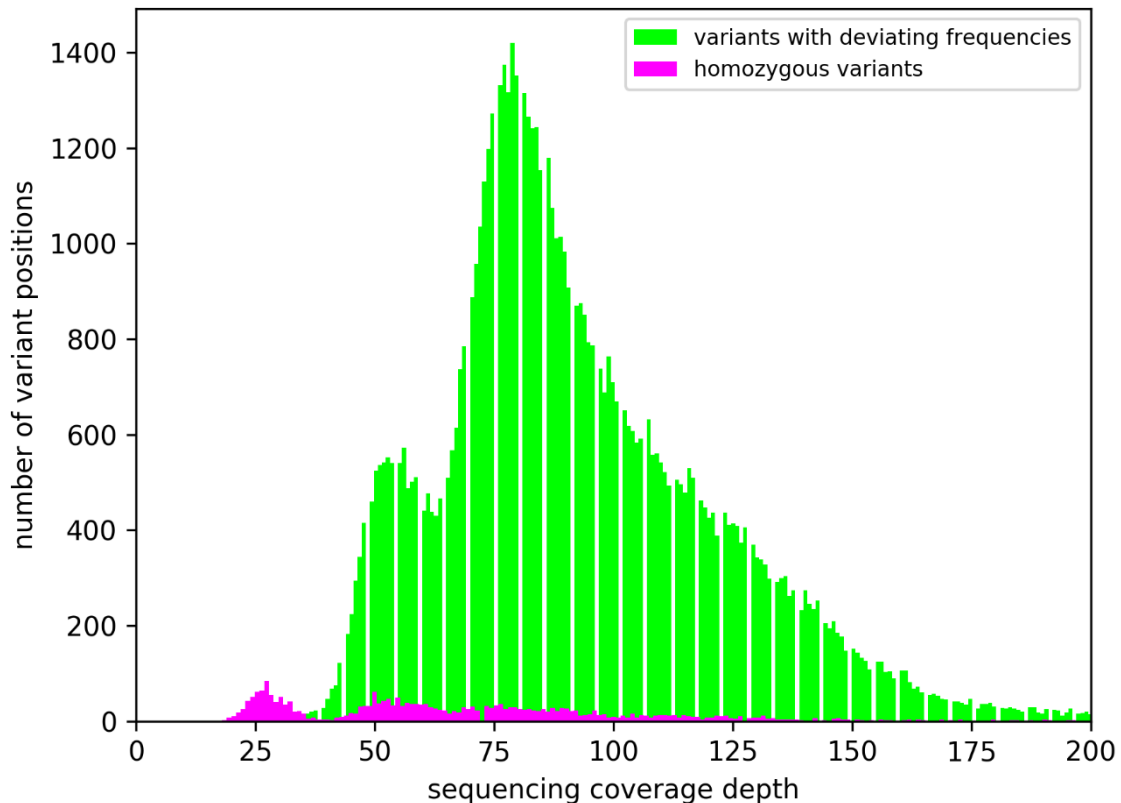
220

221

About 94.2 % of all small sequence variants display deviating frequencies. This can be explained by the asexual reproduction of the originally diploid At7 cells. A mutation only occurs in one allele and is consequently passed on to all daughter cells. However, the other allele of the same locus would still show the original sequence. It is very unlikely that the same mutation occurs in all alleles independently. Although exchange between alleles cannot be excluded, it is not a relevant mechanism for the generation of homozygous variants. Many homozygous variants between At7 and Col-0 are clustered in a few regions (Figure 2). In most cases, the respective regions are also depleted of variants with deviating frequencies. One clear example of a homozygous variant stretch is located on Chr2 between 9.5 Mbp and 11 Mbp (Figure 2, magenta marked) where the number of homozygous variants exceeds the number of variants with deviating frequencies. These 2,774 homozygous variants were probably caused by the hemizyosity of these regions in the At7 genome (Figure 1). A substantially reduced read coverage of homozygous variants compared to variants with deviating frequencies supports this hypothesis (Figure 3) and regions with high proportions of homozygous variants displayed generally low coverage. A small fraction of variants with deviating frequencies could be due to spurious read mappings and thus false positive variant calls. The weak peak of variants with about 2-fold Illumina coverage but only one detectable allele might be explained by duplication of the respective genome region after the small variant mutation occurred. The observation of a few variants with deviating frequencies in apparently hemizygous regions could be due to spontaneous mutations in some cells leading to different lineages. However, genetic drift will reduce this diversity resulting in the dominance of one allele as observed for most positions and in agreement with Muller's Ratchet [47].

A minority of the homozygous variants could also have been contributed by the original material which was used to generate this suspension cell cultures. Based on previous studies of mutation accumulation over generations [5,6,8,14], it is likely that a few of these variants have been differentiating the material used for the generation of the At7 cell culture from the material used in the *A. thaliana* genome sequencing project. Errors in

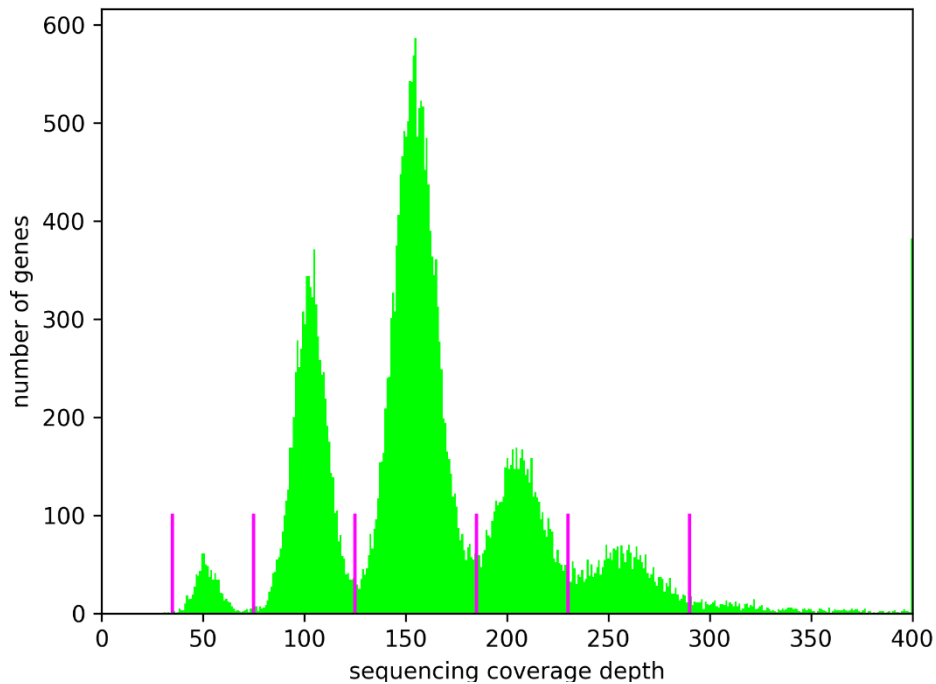
222 the At7 sequencing data and in the Col-0 reference sequence [7,11,51] could be an additional source of  
223 seemingly differences. Similar numbers were observed between genomes of individual plants of the same  
224 *A. thaliana* accessions before and have been considered to be recent events or technical artifacts [8,9].



225

226 **Figure 3. Illumina sequencing coverage depth at variant positions.** Only Illumina reads were considered in  
227 this analysis, because the variant calling was limited to this set of high quality sequences. Therefore, the average  
228 coverage of 25 fold is about half the coverage observed for the complete sequence read data set.

229 Since extremely large differences in the genome structure of At7 and Col-0 might not be revealed through  
230 variant detection tools, genes of the Col-0 genome sequence annotation Araport11 were classified based on  
231 their read mapping coverage (Figure 4, File S3 and Figure S4). About 85% of all genes display a coverage  
232 which was similar to the coverage of flanking genes. Therefore, it is likely that large genomic blocks, and not  
233 just single genes, were deleted or duplicated. This is in agreement with consistent coverage levels for large  
234 genomic blocks. Again, an average coverage of 50 fold indicates hemizygous regions, while multiples of this  
235 hemizygous coverage value indicate the presence of additional copies. The high number of genes with increased  
236 copy numbers is in agreement with previous reports of increases in chromosome duplications over cultivation  
237 time as callus [52]. Analysis of the genome of this cell line again after additional years of cultivation could  
238 reveal if there is a saturation of this increase in ploidy and nuclear DNA content. We speculate that these copy  
239 numbers are beneficial due to deleterious variants in some copies of required genes.



240

241 **Figure 4. Gene copy numbers.** Read coverage per gene was calculated as proxy for the copy number of the  
242 respective gene. Magenta lines indicate the central position of coverage valleys enclosed by peaks. Distances  
243 between these coverage valleys are not identical due to differences in peak height and resulting width  
244 differences. Absolute coverage values might be underestimated due to the removal of read pairs which appeared  
245 as the result of PCR duplicates.

### 246 3.3. Variant impact on genes and transcript isoform abundance

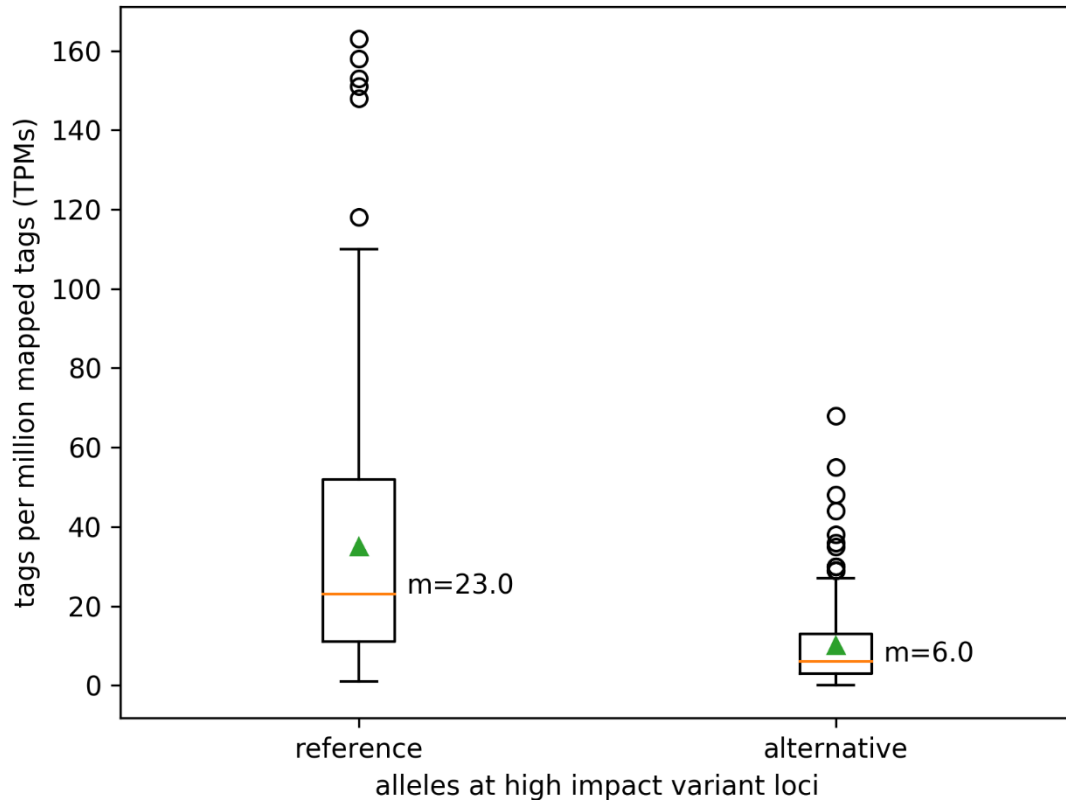
247 The impact of InDels on coding sequences was analyzed by comparing the length distributions of InDels  
248 inside protein encoding sequences with the InDel length distribution outside of these coding regions. InDel  
249 lengths which are divisible by three (one codon size) are enriched inside of coding sequences. Of all 28,213  
250 analyzed InDels, only 856 are located within the 33.5 Mbp of protein encoding regions annotated in Araport11.  
251 This depletion of InDels inside of protein encoding regions indicates that at least residual selection against  
252 disruption of these sequences is still ongoing.

253 Assessment of the functional impact of small sequence variants revealed a high impact effect (e.g.  
254 premature stop codon or frameshift) on a total of 2,189 genes (File S4). This high number can be explained by  
255 functional redundancy due to multiple alleles i.e. at least one allele is maintained in a functional state. In  
256 addition, many genes might be dispensable under stable, stress-free cell culture conditions. Therefore, the  
257 accumulation of disruptive variants or entire deletion is feasible. We restrained from gene ontology (GO)  
258 enrichment analysis due to a functional high redundancy caused by the presence of multiple alleles for most  
259 genes.

260 RNA-Seq analysis revealed a substantial difference in the abundance of native alleles and defective  
261 alleles. Isoforms with a destructive variant displayed a substantially reduced abundance (Figure 5,  
262 Mann-Whitney U-test p-value = 0). The genomic coverage of these variant positions displays a very similar  
263 pattern: Abundance of the reference allele is on average substantially higher than the coverage of the alternative  
264 alleles. This indicates that not the nonsense mediated decay pathway [53], but different copy numbers of the  
265 allele, are the main explanation for the difference in transcript abundance. Such copy number-dependent  
266 preferential allele expression is known from the vegetatively propagated autotetraploid potato [54] and seems to  
267 be a general mechanism.

268





269

270 **Figure 5. Allele specific transcript abundance at high impact variant positions.** The number of RNA-Seq  
271 reads supporting the reference (Col-0) and alternative allele, respectively, were determined at high impact  
272 variant positions.

273

#### 274 3.4 *De novo genome assembly*

275 To reveal the sequences of large structural variants, we generated an independent de novo genome  
276 assembly of At7. The Canu assembler performed best in our hands. A set of 1.4 million ONT reads longer than  
277 3 kb was assembled into 433 contigs representing 126.2 Mbp of the At7 nucleome (File S5, File S6, File S7).  
278 The N50 of 1.2 Mbp is substantially lower than other recent reports of *A. thaliana* genome assemblies [9-11].  
279 We speculate that points with abrupt changes in coverage of the At7 genome deteriorated the assembly  
280 contiguity. Manual inspection of the borders of some hemizygous regions revealed two groups of reads, which  
281 support different assembly paths: one containing the hemizygous sequence and one continuing in a different  
282 genomic location. The observation of different groups of reads is in agreement with the high number of large  
283 structural variants, which are not supported by all reads, and several distinct coverage levels detected during the  
284 read mapping to the Col-0 reference sequence. However, the total assembly size exceeds the 119.8 Mbp of the  
285 Col-0 reference sequence thus rendering the accurate anchoring of contigs in some regions, e.g. close to the  
286 centromeres, unreliable.

287

#### 288 3.5. *Future directions*

289 After the identification of genomic changes in At7 compared to the Col-0 reference genome sequence, it  
290 would be interesting to investigate genomic changes in other cell culture lines. A previous study indicates that  
291 the speed of genomic changes could be cell line specific [55]. Based on our findings and previous reports about  
292 differences in genomic stability [37], there might be certain regions in the *Arabidopsis* genome which are more  
293 likely to change than others as previously reported for *Solanum tuberosum* [56]. Aneuploidy events were  
294 described previously for an asexually reproducing plant [48] and in long term suspension and callus cultures  
295 [57,58]. Comparing genomic changes observed in independent lineages would facilitate the identification of  
296 instable or even dispensable regions. Increased copy numbers of many genomic regions seemed to have low

297 metabolic costs, but might be beneficial in buffering deleterious mutations that accumulate due to asexually  
298 reproduction.

299 In this study, we focused at the genome sequence and excluded other previously reviewed DNA  
300 modifications like methylation and chromatin condensation [59,60]. The methylation pattern is known to partly  
301 change e.g. in response to environmental conditions or developmental stage [61,62]. Changes in methylation  
302 are associated with transposable element activity [63] and general genome instability [56]. TE activation in At7  
303 could be harnessed to identify intact elements, which are repressed by methylation under normal conditions in a  
304 plant. Local hypo- and hypermethylation were reported before as the result of a tissue culture and appeared at  
305 the same genomic location across independent lines [64,65]. Generally increased methylation levels have been  
306 reported to render in vitro cultures unsuitable for long term production of secondary metabolites [66]. In the  
307 light of these reports and our genomic alteration findings, substantial methylation differences between At7 cells  
308 in a suspension culture and hypocotyl cells in a Col-0 plant can be expected. Since the number of hypocotyl  
309 cells in a plant is very limited, whole genome nanopore sequencing of this cell type to generate a reference data  
310 set for the methylation pattern is not feasible at the moment. Advances in single cell sequencing could enable  
311 such a comparison in the future.

312 In-depth investigation of genomic changes under cell cultures conditions including methylation  
313 differences could benefit applications like the development of stable transgenic lines from in vitro cell cultures  
314 [67]. Using regenerable protoplasts for non-transgenic genome editing could advance crop improvements [56]  
315 and plant cells could even be used as efficient biofactories once epigenetic challenges are overcome [66].

316

317 **Supplementary Materials: File S1:** Per chromosome histogram of At7 sequencing coverage of the Col-0 reference  
318 sequence. **File S2:** Small sequence variants between At7 and Col-0. Illumina short reads were mapped to the Col-0 reference  
319 sequence via BWA MEM. Variant calling was performed with GATK. **Figure S1:** Deletions and insertions between At7  
320 and Col-0. ONT long reads were aligned to the Col-0 reference sequence via Minimap2 and variants were identified via  
321 SVIM. **Figure S2:** Sequencing coverage of deletions between At7 and Col-0. **File S3:** Classification of Araport11 genes  
322 based on At7 read mapping coverage. **Figure S3:** Difference in sequencing coverage depth of plastome and chondrome. **File**  
323 **S4:** Selection of most severe variants based on a SnpEff prediction. **Figure S4:** Genome-wide distribution of genes with  
324 different read mapping coverage values (File S3). **File S5:** Summary of generated sequencing data. **File S6:** *De novo* At7  
325 assembly statistics. **File S7:** *De novo* At7 assembly.

326

327 **Author Contributions: Author contributions:** Conceptualization, Boas Pucker, Ralf Stracke, Jörn Kalinowski and Bernd  
328 Weisshaar; Data curation, Boas Pucker; Investigation, Boas Pucker, Christian Rückert, Ralf Stracke and Prisca Viehöver;  
329 Project administration, Boas Pucker; Software, Boas Pucker; Supervision, Bernd Weisshaar; Visualization, Boas Pucker;  
330 Writing – original draft, Boas Pucker; Writing – review & editing, Boas Pucker, Ralf Stracke and Bernd Weisshaar.

331 **Funding:** We acknowledge support for the Article Processing Charge by the Deutsche Forschungsgemeinschaft and the  
332 Open Access Publication Fund of Bielefeld University.

333 **Acknowledgments:** We thank Daniela Holtgräwe and Hanna Marie Schilbert for great support.

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

## 335 References

- 336 1. The Arabidopsis Genome Initiative. Analysis of the genome sequence of the flowering plant  
337 *Arabidopsis thaliana*. *Nature* **2000**, *408*, 796-815.
- 338 2. Schneeberger, K.; Ossowski, S.; Ott, F.; Klein, J.D.; Wang, X.; Lanz, C.; Smith, L.M.; Cao, J.; Fitz, J.;  
339 Warthmann, N., et al. Reference-guided assembly of four diverse *Arabidopsis thaliana* genomes.  
340 *Proceedings of the National Academies of Sciences of the United States of America* **2011**, *108*,  
341 10249-10254.
- 342 3. Consortium, T.G. 1,135 Genomes Reveal the Global Pattern of Polymorphism in *Arabidopsis thaliana*.  
343 *Cell* **2016**, *166*, 481-491.
- 344 4. Durvasula, A.; Fulgione, A.; Gutaker, R.M.; Alacakaptan, S.I.; Flood, P.J.; Neto, C.; Tsuchimatsu, T.;  
345 Burbano, H.A.; Picó, F.X.; Alonso-Blanco, C., et al. African genomes illuminate the early history and

- 346 transition to selfing in *Arabidopsis thaliana*. *Proceedings of the National Academy of Sciences of the*  
347 *United States of America* **2017**, *114*, 5213-5218.
- 348 5. Ossowski, S.; Schneeberger, K.; Clark, R.; Lanz, C.; Warthmann, N.; Weigel, D. Sequencing of  
349 natural strains of *Arabidopsis thaliana* with short reads. *Genome Research* **2008**, *18*, 2024–2033.
- 350 6. Cao, J.; Schneeberger, K.; Ossowski, S.; Günther, T.; Bender, S.; Fitz, J.; Koenig, D.; Lanz, C.; Stegle,  
351 O.; Lippert, C., et al. Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. *Nature*  
352 *Genetics* **2011**, *43*, 956-963.
- 353 7. Long, Q.; Rabanal, F.A.; Meng, D.; Huber, C.D.; Farlow, A.; Platzer, A.; Zhang, Q.; Vilhjálmsson,  
354 B.J.; Korte, A.; Nizhynska, V., et al. Massive genomic variation and strong selection in *Arabidopsis*  
355 *thaliana* lines from Sweden. *Nature Genetics* **2013**, *45*, 884-890.
- 356 8. Pucker, B.; Holtgräwe, D.; Rosleff Sörensen, T.; Stracke, R.; Viehöver, P.; Weisshaar, B. A De Novo  
357 Genome Sequence Assembly of the *Arabidopsis thaliana* Accession Niederzenz-1 Displays  
358 Presence/Absence Variation and Strong Synteny. *PLoS ONE* **2016**, *11*, e0164321.
- 359 9. Zapata, L.; Ding, J.; Willing, E.M.; Hartwig, B.; Bezdán, D.; Jiao, W.B.; Patel, V.; Velikkakam James,  
360 G.; Koornneef, M.; Ossowski, S., et al. Chromosome-level assembly of *Arabidopsis thaliana* Ler  
361 reveals the extent of translocation and inversion polymorphisms. *Proceedings of the National*  
362 *Academy of Sciences of the United States of America* **2016**, *113*, 4052-4060.
- 363 10. Michael, T.P.; Jupe, F.; Bemm, F.; Motley, S.T.; Sandoval, J.P.; Lanz, C.; Loudet, O.; Weigel, D.;  
364 Ecker, J.R. High contiguity *Arabidopsis thaliana* genome assembly with a single nanopore flow cell.  
365 *Nature Communications* **2018**, *9*, 541.
- 366 11. Pucker, B.; Holtgräwe, D.; Stadermann, K.B.; Frey, K.; Huettel, B.; Reinhardt, R.; Weisshaar, B. A  
367 chromosome-level sequence assembly reveals the structure of the *Arabidopsis thaliana* Nd-1 genome  
368 and its gene set. *PLoS One* **2019**, *14*, e0216233.
- 369 12. Debladis, E.; Llauro, C.; Carpentier, M.C.; Mirouze, M.; Panaud, O. Detection of active transposable  
370 elements in *Arabidopsis thaliana* using Oxford Nanopore Sequencing technology. *BMC Genomics*  
371 **2017**, *18*, 537.
- 372 13. Ossowski, S.; Schneeberger, K.; Lucas-Lledó, J.I.; Warthmann, N.; Clark, R.M.; Shaw, R.G.; Weigel,  
373 D.; Lynch, M. The rate and molecular spectrum of spontaneous mutations in *Arabidopsis thaliana*.  
374 *Science* **2010**, *327*, 92-94.
- 375 14. Exposito-Alonso, M.; Becker, C.; Schuenemann, V.J.; Reiter, E.; Setzer, C.; Slovak, R.; Brachi, B.;  
376 Hagemann, J.; Grimm, D.G.; Chen, J., et al. The rate and potential relevance of new mutations in a  
377 colonizing plant lineage. *PLoS Genetics* **2018**, *14*, e1007155.
- 378 15. Dangl, J.L.; Hauffe, K.-D.; Lipphardt, S.; Hahlbrock, K.; Scheel, D. Parsley protoplasts retain  
379 differential responsiveness to UV light and fungal elicitor. *The EMBO Journal* **1987**, *6*, 2551-2556.
- 380 16. Hartmann, U.; Valentine, W.J.; Christie, J.M.; Hays, J.; Jenkins, G.I.; Weisshaar, B. Identification of  
381 UV/blue light-response elements in the *Arabidopsis thaliana* chalcone synthase promoter using a  
382 homologous protoplast transient expression system. *Plant Molecular Biology* **1998**, *36*, 741-754.
- 383 17. Baudry, A.; Heim, M.A.; Dubreucq, B.; Caboche, M.; Weisshaar, B.; Lepiniec, L. TT2, TT8, and  
384 TTG1 synergistically specify the expression of *BANYULS* and proanthocyanidin biosynthesis in  
385 *Arabidopsis thaliana*. *The Plant Journal* **2004**, *39*, 366-380.
- 386 18. Trezzini, G.F.; Horrichs, A.; Somssich, I.E. Isolation of putative defense-related genes from  
387 *Arabidopsis thaliana* and expression in fungal elicitor-treated cells. *Plant Molecular Biology* **1993**, *21*,  
388 385-389.

- 389 19. Stracke, R.; Thiedig, K.; Kuhlmann, M.; Weisshaar, B. Analyzing synthetic promoters using  
390 Arabidopsis protoplasts. In *Methods in Molecular Biology: Plant Synthetic Promoters*, 2016/08/26  
391 ed.; Hehl, R., Ed. Springer: New York, 2016; Vol. 1482, pp. 67-81.
- 392 20. Rosso, M.G.; Li, Y.; Strizhov, N.; Reiss, B.; Dekker, K.; Weisshaar, B. An *Arabidopsis thaliana*  
393 T-DNA mutagenised population (GABI-Kat) for flanking sequence tag based reverse genetics. *Plant*  
394 *Molecular Biology* **2003**, *53*, 247-259.
- 395 21. Koren, S.; Walenz, B.P.; Berlin, K.; Miller, J.R.; Bergman, N.H.; Phillippy, A.M. Canu: scalable and  
396 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research*  
397 **2017**, *27*, 722-736.
- 398 22. Li, H. Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences.  
399 *Bioinformatics* **2016**, *32*, 2103-2110.
- 400 23. Kolmogorov, M.; Yuan, J.; Lin, Y.R.; Pevzner, P.A. Assembly of Long Error-Prone Reads Using  
401 Repeat Graphs. *Nature Biotechnology* **2019**, *37*, 540-546.
- 402 24. Loman, N.J.; Quick, J.; Simpson, J.T. A complete bacterial genome assembled de novo using only  
403 nanopore sequencing data. *Nature Methods* **2015**, *12*, 733-735.
- 404 25. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *Oxford*  
405 *University Press* **2013**, 1-3.
- 406 26. Walker, B.J.; Abeel, T.; Shea, T.; Priest, M.; Abouelliel, A.; Sakthikumar, S.; Cuomo, C.A.; Zeng, Q.;  
407 Wortman, J.; Young, S.K., et al. Pilon: an integrated tool for comprehensive microbial variant  
408 detection and genome assembly improvement. *PLoS ONE* **2014**, *9*, e112963.
- 409 27. McKenna, A.; Hanna, M.; Banks, E.; Sivachenko, A.; Cibulskis, K.; Kernysky, A.; Garimella, K.;  
410 Altshuler, D.; Gabriel, S.; Daly, M., et al. The Genome Analysis Toolkit: a MapReduce framework for  
411 analyzing next-generation DNA sequencing data. *Genome Research* **2010**, *20*, 1297-1303.
- 412 28. Van der Auwera, G.A.; Carneiro, M.O.; Hartl, C.; Poplin, R.; Del Angel, G.; Levy-Moonshine, A.;  
413 Jordan, T.; Shakir, K.; Roazen, D.; Thibault, J., et al. From FastQ data to high confidence variant calls:  
414 the Genome Analysis Toolkit best practices pipeline. *Current Protocols in Bioinformatics* **2013**, *11*,  
415 1110.
- 416 29. Baasner, J.S.; Howard, D.; Pucker, B. Influence of neighboring small sequence variants on functional  
417 impact prediction. *bioRxiv* **2019**, <http://dx.doi.org/10.1101/596718>.
- 418 30. Cingolani, P.; Platts, A.; Wang le, L.; Coon, M.; Nguyen, T.; Wang, L.; Land, S.J.; Lu, X.; Ruden,  
419 D.M. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:  
420 SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* **2012**, *6*,  
421 80-92.
- 422 31. Heller, D.; Vingron, M. SVIM: Structural Variant Identification using Mapped Long Reads.  
423 *Bioinformatics* **2019**, btz041.
- 424 32. Dobin, A.; Davis, C.A.; Schlesinger, F.; Drenkow, J.; Zaleski, C.; Jha, S.; Batut, P.; Chaisson, M.;  
425 Gingeras, T.R. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **2013**, *29*, 15-21.
- 426 33. Haak, M.; Vinke, S.; Keller, W.; Droste, J.; Rückert, C.; Kalinowski, J.; Pucker, B. High Quality de  
427 Novo Transcriptome Assembly of *Croton tiglium*. *Frontiers in Molecular Biosciences* **2018**, *5*, 62.
- 428 34. Li, H.; Handsaker, B.; Wysoker, A.; Fennell, T.; Ruan, J.; Homer, N.; Marth, G.; Abecasis, G.; Durbin,  
429 R. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **2009**, *25*, 2078-2079.
- 430 35. Quinlan, A.R.; Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic features.  
431 *Bioinformatics* **2010**, *26*, 841-842.

- 432 36. Lee, H.; McManus, C.J.; Cho, D.Y.; Eaton, M.; Renda, F.; Somma, M.P.; Cherbas, L.; May, G.;  
433 Powell, S.; Zhang, D., et al. DNA copy number evolution in *Drosophila* cell lines. *Genome Biology*  
434 **2014**, *15*, R70.
- 435 37. Kaas, C.S.; Kristensen, C.; Betenbaugh, M.J.; Andersen, M.R. Sequencing the CHO DXB11 genome  
436 reveals regional variations in genomic stability and haploidy. *BMC Genomics* **2015**, *16*, 160.
- 437 38. DAmato, F. Nuclear changes in cultured plant cells. *Caryologia* **1991**, *44*, 217-224.
- 438 39. Tan, E.H.; Henry, I.M.; Ravi, M.; Bradnam, K.R.; Mandakova, T.; Marimuthu, M.P.; Korf, I.; Lysak,  
439 M.A.; Comai, L.; Chan, S.W. Catastrophic chromosomal restructuring during genome elimination in  
440 plants. *Elife* **2015**, *4*, doi:10.7554/eLife.06516.
- 441 40. Lynch, M.; Marinov, G.K. The bioenergetic costs of a gene. *Proceedings of the National Academy of*  
442 *Sciences of the United States of America* **2015**, *112*, 15690-15695.
- 443 41. Hu, T.T.; Pattyn, P.; Bakker, E.G.; Cao, J.; Cheng, J.F.; Clark, R.M.; Fahlgren, N.; Fawcett, J.A.;  
444 Grimwood, J.; Gundlach, H., et al. The *Arabidopsis lyrata* genome sequence and the basis of rapid  
445 genome size change. *Nature Genetics* **2011**, *43*, 476-481.
- 446 42. Ohlrogge, J.B.; Kuhn, D.N.; Stumpf, P.K. Subcellular localization of acyl carrier protein in leaf  
447 protoplasts of *Spinacia oleracea*. *Proceedings of the National Academy of Sciences of the United States*  
448 *of America* **1979**, *76*, 1194-1198.
- 449 43. Rawsthorne, S. Carbon flux and fatty acid synthesis in plants. *Progress in Lipid Research* **2002**, *41*,  
450 182-196.
- 451 44. Lichtenthaler, H.K.; Schwender, J.; Disch, A.; Rohmer, M. Biosynthesis of isoprenoids in higher plant  
452 chloroplasts proceeds via a mevalonate-independent pathway. *FEBS Letters* **1997**, *400*, 271-274.
- 453 45. Cole, L.W. The Evolution of Per-cell Organelle Number. *Frontiers in Cell and Developmental Biology*  
454 **2016**, *4*.
- 455 46. Kleinboelting, N.; Huep, G.; Appelhagen, I.; Viehoveer, P.; Li, Y.; Weisshaar, B. The Structural  
456 Features of Thousands of T-DNA Insertion Sites Are Consistent with a Double-Strand Break  
457 Repair-Based Insertion Mechanism. *Molecular Plant* **2015**, *8*, 1651-1664.
- 458 47. Muller, H.J. Some Genetic Aspects of Sex. *The American Naturalist* **1932**, *66*, 118-138.
- 459 48. Schranz, M.E.; Kantama, L.; de Jong, H.; Mitchell-Olds, T. Asexual reproduction in a close relative of  
460 *Arabidopsis*: a genetic investigation of apomixis in *Boechera* (Brassicaceae). *New Phytologist* **2006**,  
461 *171*, 425-438.
- 462 49. Lovell, J.T.; Williamson, R.J.; Wright, S.I.; McKay, J.K.; Sharbel, T.F. Mutation Accumulation in an  
463 Asexual Relative of *Arabidopsis*. *PLoS Genetics* **2017**, *13*, e1006550.
- 464 50. Bomblies, K.; Yant, L.; Laitinen, R.A.; Kim, S.T.; Hollister, J.D.; Warthmann, N.; Fitz, J.; Weigel, D.  
465 Local-Scale Patterns of Genetic Variability, Outcrossing, and Spatial Structure in Natural Stands of  
466 *Arabidopsis thaliana*. *PLoS Genetics* **2010**, *6*, e1000890.
- 467 51. Kawakatsu, T.; Huang, S.S.; Jupe, F.; Sasaki, E.; Schmitz, R.J.; Urich, M.A.; Castanon, R.; Nery, J.R.;  
468 Barragan, C.; He, Y., et al. Epigenomic Diversity in a Global Collection of *Arabidopsis thaliana*  
469 Accessions. *Cell* **2016**, *166*, 492-505.
- 470 52. Molina, M.; Garcia, M.D. Analysis of Genetic Variability and Regenerated in Long-term Callus  
471 Cultures Plants of Maize. *Cytologia* **1998**, *63*, 183-190.
- 472 53. Hug, N.; Longman, D.; Cáceres, J.F. Mechanism and regulation of the nonsense-mediated decay  
473 pathway. *Nucleic acids research* **2016**, *44*, 1483-1495.

- 474 54. Pham, G.M.; Newton, L.; Wiegert-Rininger, K.; Vaillancourt, B.; Douches, D.S.; Buell, C.R.  
475 Extensive genome heterogeneity leads to preferential allele expression and copy number-dependent  
476 expression in cultivated potato. *The Plant Journal* **2017**, *92*, 624-637.
- 477 55. Zucchi, M.I.; Arizono, H.; Morais, V.A.; Pelegri-nelli Fungaro, M.H.; Carneiro Vieira, M.L. Genetic  
478 instability of sugarcane plants derived from meristem cultures. *Genetics and Molecular Biology* **2002**,  
479 *25*, 91-96.
- 480 56. Fossi, M.; Amundson, K.; Kuppu, S.; Britt, A.; Comai, L. Regeneration of *Solanum tuberosum* Plants  
481 from Protoplasts Induces Widespread Genome Instability. *Plant Physiology* **2019**, *180*, 78-86.
- 482 57. Giorgetti, L.; Ruffini Castiglione, M.; Turrini, A.; Ronchi, V.N.; Geri, C. Cytogenetic and histological  
483 approach for early detection of “mantled” somaclonal variants of oil palm regenerated by somatic  
484 embryogenesis: first results on the characterization of regeneration system. *Caryologia* **2011**, *64*,  
485 223-234.
- 486 58. Landey, R.B.; Cenci, A.; Guyot, R.; Bertrand, B.; Georget, F.; Dechamp, E.; J.-C., H.; Aribi, J.;  
487 Lashermes, P.; Etienne, H. Assessment of genetic and epigenetic changes during cell culture ageing  
488 and relations with somaclonal variation in *Coffea arabica*. *Plant Cell Tiss Organ Cult* **2015**, *122*,  
489 517–531.
- 490 59. Miguel, C.; Marum, L. An epigenetic view of plant cells cultured in vitro: somaclonal variation and  
491 beyond. *Journal of Experimental Botany* **2011**, *62*, 3713-3725.
- 492 60. Neelakandan, A.K.; Wang, K. Recent progress in the understanding of tissue culture-induced genome  
493 level changes in plants and potential applications. *Plant Cell Rep* **2012**, *31*, 597-620.
- 494 61. Bartels, A.; Han, Q.; Nair, P.; Stacey, L.; Gaynier, H.; Mosley, M.; Huang, Q.Q.; Pearson, J.K.; Hsieh,  
495 T.F.; An, Y.C., et al. Dynamic DNA Methylation in Plant Growth and Development. *International*  
496 *Journal of Molecular Sciences* **2018**, *19*.
- 497 62. Thiebaut, F.; Hemerly, A.S.; Ferreira, P.C.G. A Role for Epigenetic Regulation in the Adaptation and  
498 Stress Responses of Non-model Plants. *Frontiers in Plant Science* **2019**, *10*, 246.
- 499 63. Springer, N.M.; Lisch, D.; Li, Q. Creating Order from Chaos: Epigenome Dynamics in Plants with  
500 Complex Genomes. *The Plant Cell* **2016**, *28*, 314-325.
- 501 64. Stelpflug, S.C.; Eichten, S.R.; Hermanson, P.J.; Springer, N.M.; Kaeppler, S.M. Consistent and  
502 heritable alterations of DNA methylation are induced by tissue culture in maize. *Genetics* **2014**, *198*,  
503 209-218.
- 504 65. Han, Z.; Crisp, P.A.; Stelpflug, S.; Kaeppler, S.M.; Li, Q.; Springer, N.M. Heritable Epigenomic  
505 Changes to the Maize Methylome Resulting from Tissue Culture. *Genetics* **2018**, *209*, 983-995.
- 506 66. Sanchez-Munoz, R.; Moyano, E.; Khojasteh, A.; Bonfill, M.; Cusido, R.M.; Palazon, J. Genomic  
507 methylation in plant cell cultures: A barrier to the development of commercial long-term biofactories.  
508 *Eng Life Sci* **2019**, *1*, 8.
- 509 67. Stroud, H.; Ding, B.; Simon, S.A.; Feng, S.; Bellizzi, M.; Pellegrini, M.; Wang, G.L.; Meyers, B.C.;  
510 Jacobsen, S.E. Plants regenerated from tissue culture contain stable epigenome changes in rice. *eLife*  
511 **2013**, *2*, e00354.