# 1 Extensive gene duplication in
# 2 Arabidopsis revealed by pseudo-
# 3 heterozygosity

4 Benjamin Jaegle [1], Luz Mayela Soto-Jiménez [1], Robin Burns [1][3], Fernando A. Rabanal [2],
5 Magnus Nordborg [1]

# 6 Abstract

7 **Background:** It is becoming apparent that genomes harbor massive amounts of structural
8 variation, and that this variation has largely gone undetected for technical reasons. In addition to
9 being inherently interesting, structural variation can cause artifacts when short-read sequencing
10 data are mapped to a reference genome. In particular, spurious SNPs (that do not show
11 Mendelian segregation) may result from mapping of reads to duplicated regions. Recalling SNP
12 using the raw reads of the 1001 Arabidopsis Genomes Project we identified 3.3 million
13 heterozygous SNPs (44% of total). Given that *Arabidopsis thaliana* (*A. thaliana*) is highly selfing,
14 we hypothesized that these SNPs reflected cryptic copy number variation, and investigated
15 them further.

16 **Results:** While genuine heterozygosity should occur in tracts within individuals, heterozygosity
17 at a particular locus is instead shared across individuals in a manner that strongly suggests it
18 reflects segregating duplications rather than actual heterozygosity. Focusing on pseudo-
19 heterozygosity in annotated genes, we used GWAS to map the position of the duplicates,
20 identifying 2500 putatively duplicated genes. The results were validated using *de novo* genome
21 assemblies from six lines. Specific examples included an annotated gene and nearby
22 transposon that, in fact, transpose together.

23 **Conclusions:** Our study confirms that most heterozygous SNPs calls in *A. thaliana* are
24 artifacts, and suggest that great caution is needed when analysing SNP data from short-read
25 sequencing. The finding that 10% of annotated genes are copy-number variables, and the

26 realization that neither gene- nor transposon-annotation necessarily tells us what is actually
27 mobile in the genome suggest that future analyses based on independently assembled
28 genomes will be very informative.

29 **Keywords:** structural variation, gene duplication, GWAS, SNP calling

# 30 Introduction

31 With the sequencing of genomes becoming routine, it is evident that, besides single nucleotide
32 polymorphisms (SNPs), structural variants (SVs) play a major role in genome variation (Alkan,
33 Coe, and Eichler 2011). There are many kinds of SVs, e.g., indels, inversions, and
34 transpositions. Of particular interest from a functional point of view is gene duplication, leading
35 to copy number variation (CNV).

36 Before Next-Generation Sequencing (NGS) was available, genome-wide detection of
37 CNVs was achieved using DNA-microarrays. These methods had severe weaknesses, leading
38 to low resolution and problems detecting novel and rare mutations. (Carter 2007; Snijders et al.
39 2001). With the development of NGS, our ability to detect CNVs increased dramatically, using
40 tools based on split reads, sequencing coverage, or even *de novo* assembly (Shendure and Ji
41 2008; Zhao et al. 2013). In mammals, many examples of CNVs with a major phenotypic effect
42 have been found (Gonzalez et al. 2005; Perry et al. 2007; Handsaker et al. 2011). One example
43 is the duplication of MWS/MLS associated with better trichromatic color vision (Miyahara et al.
44 1998).

45 While early investigation of CNV focused on mammals, several subsequent studies have
46 looked at plant genomes. In *Brassica rapa,* gene CNV has been shown to be involved in
47 morphological variation (Lin et al. 2014) and an analysis of the poplar pan genome revealed at
48 least 3000 genes affected by CNV (Pinosio et al. 2016). It has also been shown that variable
49 regions in the rice genome are enriched in genes related to defence to biotic stress. (Yao et al.
50 2015). More recently, the first chromosome-level assemblies of seven accessions of *A. thaliana*
51 based on long-read sequencing were released (Jiao and Schneeberger 2019), demonstrating
52 that a large proportion of the genome is structurally variable. Similar studies have also been
53 carried out in maize (C. Li et al. 2020; Hufford et al. 2021), tomato (Alonge et al. 2020), rice
54 (Zhou et al. 2020) and soybean (Y. Liu et al. 2020). These approaches are likely to provide a
55 more comprehensive picture than short-read sequencing, but are also far more expensive.

56    In 2016, the 1001 Genomes Consortium released short-read sequencing data and SNP
57    calls for 1135 *A. thaliana* accessions (1001 Genomes Consortium 2016). Several groups have
58    used these data to identify large numbers of structural variants using split reads (Göktay,
59    Fulgione, and Hancock 2020; Zmienko et al. 2020; D.-X. Liu et al. 2021). Here we approach this
60    from a different angle. Our starting point is the startling observation that recalling SNPs using
61    the raw reads of the 1001 Genomes data set we identified 3.3 million (44% of total) putatively
62    heterozygous SNPs. In a highly selfing organism, this is obviously highly implausible, and these
63    SNPs were flagged as spurious, presumably products of cryptic CNV, which can generate
64    "pseudo-SNPs" (Ranade et al. 2001) when sequencing reads from non-identical duplicates are
65    (mis-)mapped to a reference genome that does not contain the duplication. Note that allelic SNP
66    differences are expected to exist *ab initio* in the population, leading to instant pseudo-
67    heterozygosity as soon as the duplicated copy recombines away from its template. In this paper
68    we return to these putative pseudo-SNPs and show that they are indeed largely due to
69    duplications, the position of which can be precisely mapped using GWAS. Our approach is
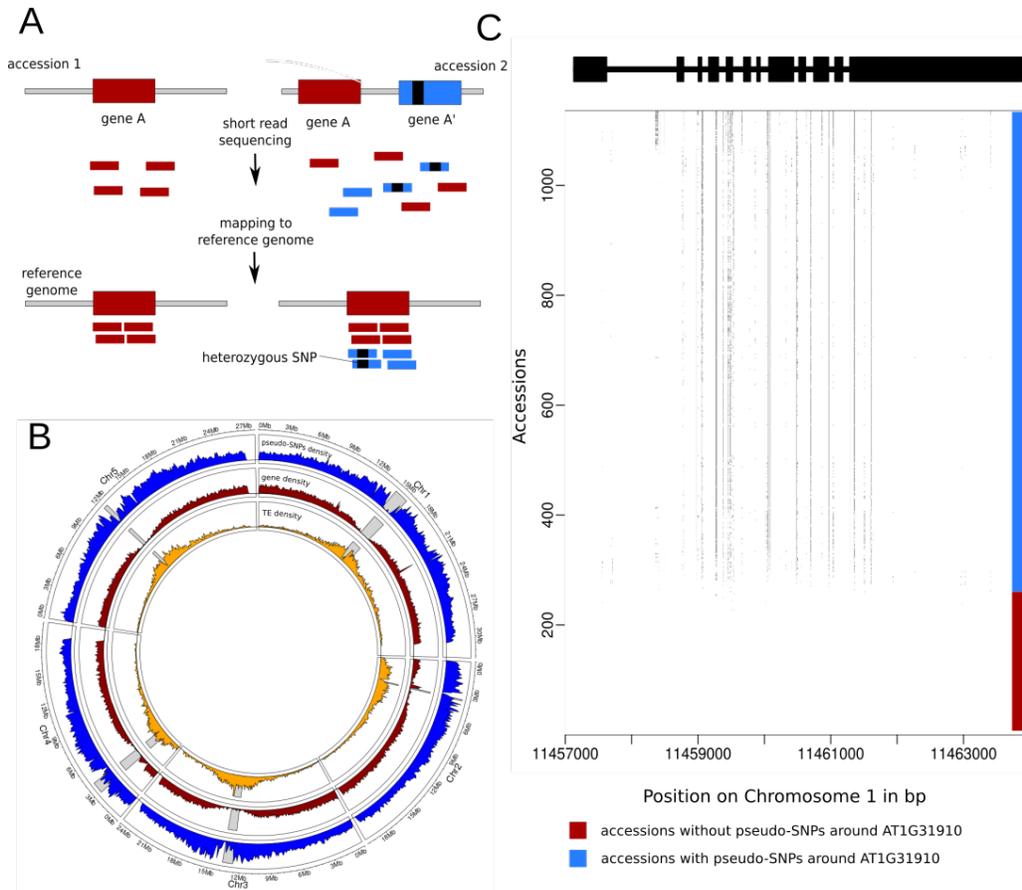70    broadly applicable, and we demonstrate that it can reveal interesting biology.

# 71    Analysis

## 72    Massive pseudo-heterozygosity in the 1001 Genomes data

73    Given that *A. thaliana* is highly selfing, a large fraction (44%) of heterozygous SNPs is
74    inherently implausible. Two other lines of evidence support the conclusion that they are
75    spurious. First, genuine residual heterozygosity would appear as large genomic tracts of
76    heterozygosity in individuals with recent outcrossing in their ancestry. Being simply a product of
77    recombination and Mendelian segregation, which tracts remain heterozygous is random, and
78    there is no reason two individuals would share tracts unless they are very closely related. The
79    observed pattern is completely the opposite. While a small number of individuals do show signs
80    of recent outcrossing, this is quite rare (as expected given the low rate of outcrossing in this
81    species, and the fact that the sequenced individuals were selected to be completely inbred).

82    Instead we find that the same SNP are often heterozygous in multiple individuals.
83    Although the population frequency of heterozygosity at a given SNP is typically low
84    **(Supplemental Figure 1)**, over a million heterozygous SNPs are shared by at least 5
85    accessions, and a closer look at the pattern of putative heterozygosity usually reveals short

84    tracts of shared heterozygosity that would be vanishingly unlikely under residual heterozygosity,

85    but would be expected if the tract represents a shared duplication, and heterozygosity is in fact

86    pseudo-heterozygosity due to mis-mapped reads (**Figure 1**).



87    **Figure 1:** Pseudo-heterozygosity in the 1001 Genomes dataset. (**A**) Cartoon illustrating how a duplication

88    can generate pseudo-SNPs when mapping to a reference genome that does not contain the duplication.

89    (**B**) Genomic density of transposons, genes, and shared heterozygous SNPs. (**C**) The pattern of putative

90    heterozygosity around AT1G31910 for the 1057 accessions. Dots in the plot represent putative

91    heterozygosity.

92         Furthermore, the density of shared heterozygous SNPs is considerably higher around

93    the centromeres (**Figure 1**), which is again not expected under random residual heterozygosity,

94    but is rather reminiscent of the pattern observed for transposons, where it is interpreted as the

95    result of selection removing insertions from euchromatic regions, leading to a build-up of

96    common (shared) transposon insertions near centromere (Quadrana et al. 2016). As we shall
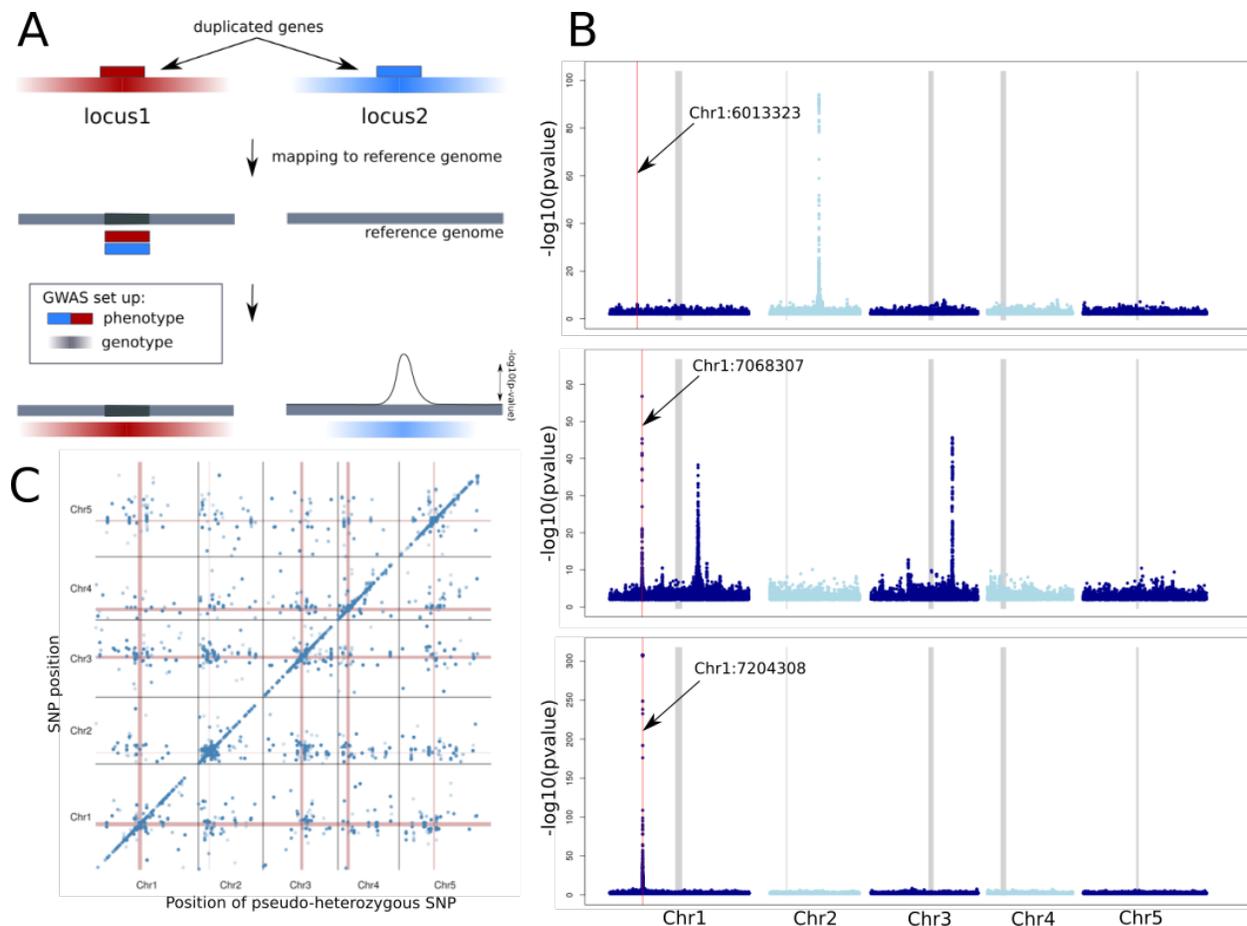
97    see below, it is likely that transposons play an important role in generating cryptic duplications

98    leading to pseudo-heterozygosity.

99        Despite the evidence for selection against these putative duplications, we found 2570

100   genes containing pseudo-SNPs segregating at 5% or more in the population (**Supplemental**

101   **Figure 2**). Gene-ontology analysis of these genes reveals an enrichment for biological

102   processes involved in response to UV-B, bacteria or fungi (**Supplemental Figure 3**). In the

103   following sections, we investigate these putatively duplicated genes further.

# Mapping common duplications using genome-wide association

105   If heterozygosity is caused by the presence of cryptic duplications in non-reference genomes, it

106   should be possible to map the latter using GWAS and heterozygosity as a "phenotype"

107   (Imprialou et al 2017). We did this for each of a total of 26647 SNPs exhibiting heterozygosity

108   within the genes described above.

109        Of the 2570 genes that showed evidence of duplication, 2511 contained at least one

110   major association (using significance threshold of $p < 10^{-20}$; see Methods). For 708 genes, the

111   association was more than 50 kb away from the pseudo-SNP used to define the phenotype, and

112   for 175 it was within 50 kb. We will refer to these as *trans-* and *cis-*associations, respectively.

113   The majority of genes, 1628, had both *cis-* and *trans-*associations (**Figure 2**).

**Figure 2:** GWAS of putative duplications A: Schematic representation of the principle of how GWAS can be used to detect the position of the duplicated genes based on linkage disequilibrium (LD). As phenotype, heterozygosity at the position of interest is coded as 1 (present) or 0 (absent). As genotype, the SNPs matrix of the 1001 genome dataset was used (with heterozygous SNPs filtered out). Color gradients represent the strength of LD around the two loci. In this example the reference genome does not contain locus2. (**B**) GWAS results for three different genes with evidence of duplication. The grey boxes represent the centromere of each chromosome. The red lines indicate the position of the pseudo-SNP used for each GWAS and the thick grey lines indicate the centromeres. The top plot shows a *trans*-association, the bottom a *cis*-association, and the middle shows a case with both (*cis* plus two *trans*). The precise coordinates in base-pair on chromosome 1 are 6013323, 7068307 and 7204308. (**C**) All 26647 GWAS results summary.

To validate these results we assembled 6 non-reference genomes *de novo* using long-read PacBio sequencing. The GWAS hits tells us where we should expect to find the duplication (the cause of pseudo-heterozygosity) using coordinates from the reference genome. We identified the homologous region of each non-reference genome, then used BLAST to search

129 for evidence of duplication. Of the 398 genes predicted to have a duplication present in at least
130 one of the 6 non-reference genomes, 235 (60%) were found to do so. The distribution of
131 fragment sizes detected suggests that we capture a mixture of gene fragments and full genes
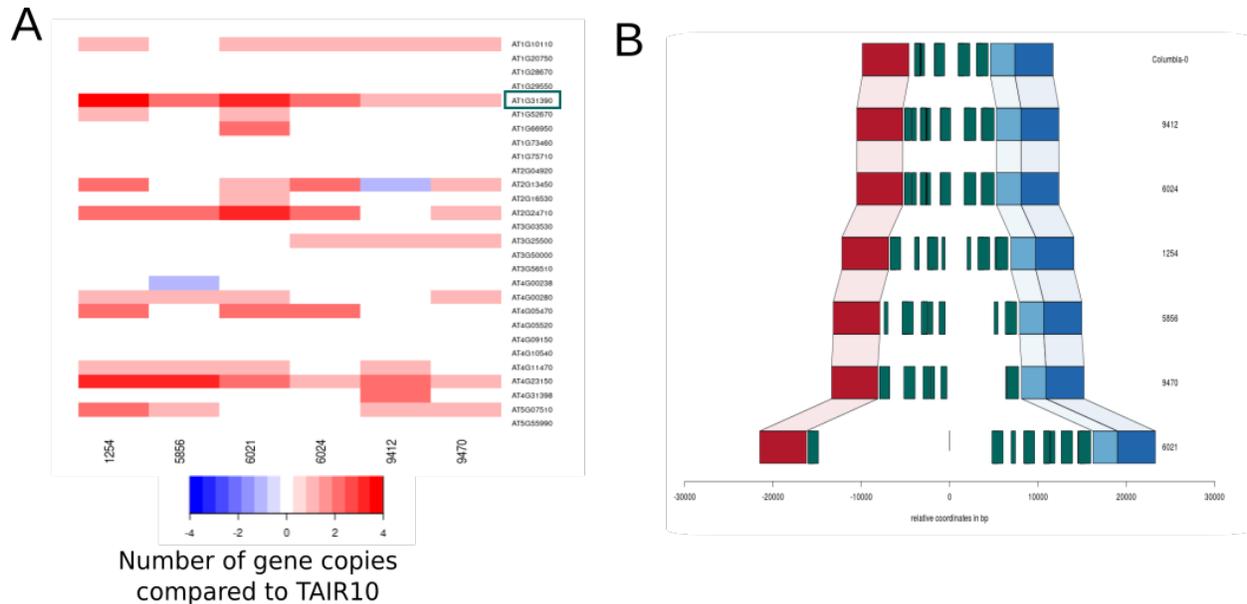132 (**Supplemental Figure 4**).

## 133 Rare duplications

134 The GWAS approach has no power to detect rare duplications, which is why we restricted the
135 analysis above to pseudo-heterozygous SNPs seen in five or more individuals. Yet most are
136 rarer: 40% are seen only in a single individual, and 16% are seen in two. As it turns out, many of
137 these appear to be associated with more common duplications. Restricting ourselves to genes
138 only, 11.4% of the singleton pseudo-heterozygous SNPs are found in the 2570 genes already
139 identified using common duplications, a significant excess ($p$ = 2.481877e-109). For doubletons,
140 the percentage is 11.1% ($p$ = 1.882515e-139). Whether they are caused by the same
141 duplications, or reflect additional ones present at lower frequency is impossible to say. To
142 confirm duplications more directly, we simply took the reads generating the singleton and
143 doubleton pseudo-heterozygotes, and compared the result of mapping them to the reference
144 genome, and to the right genome (from the same inbred line).  A consequence of the reads
145 mapping at different locations is that the mapping coverage around the pseudo-SNPs will be
146 decreased when mapping to the PacBio genomes. As expected a high proportion of the SNPs
147 tested have lower coverage when mapping to the PacBio genomes **(Supplemental Figure 5-6)**.
148 On top of the decrease in coverage we could also detect reads mapping to multiple locations as
149 well as the disappearance of the Pseudo-SNPs. Overall, 41.5% of the doubletons tag regions
150 that map in more regions in the PacBio genomes compared to the reference genome
151 (**Supplemental Figure 5-6, 7**).

## 152 Local duplications

153 If duplications arise via tandem duplications, they will not give rise to pseudo-SNPs until the
154 copies have diverged via mutations. This is in contrast to unlinked copies, which will lead to
155 pseudo-SNPs due to standing allelic variation as soon as recombination has separated copy
156 from original. We should thus expect the approach taken here to be biased against detecting
157 local duplications. Nonetheless, GWAS revealed 175 genes with evidence only for a *cis*
158 duplication. 28 of these were predicted to be present in at least one of the 6 PacBio assemblies,
159 and 16 could be confirmed to have local variation of copy number compare to the reference.

160    (**Figure 4A**).



**Figure 3:** Confirmation of tandem duplications. (**A**) The distribution of estimated copy number (based on sequencing coverage) across 6 PacBio genomes for 28 genes predicted to be involved in tandem duplications based on the analyses of this paper. (**B**) The duplication pattern observed in these genomes for one such gene, AT1G31390. Each color represents gene annotation from TAIR10. Multiple rectangles of the same copies denote multiple copies of the same gene.

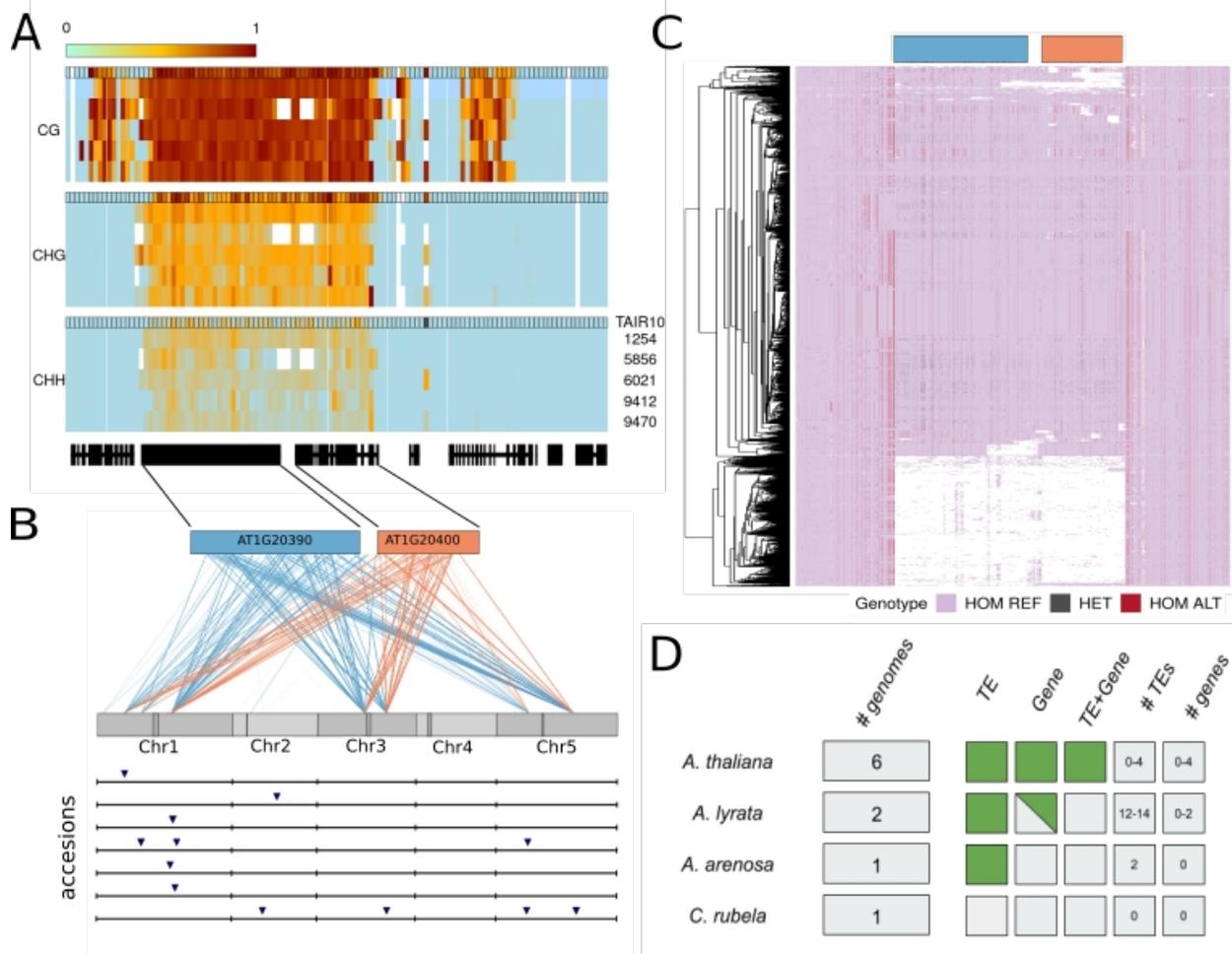166        The local structure of the duplications can be complex. An example is provided by the gene AT1G31390, annotated as a member of MATH/TRAF-domain genes, and which appears to be present in 4 tandem copies in the reference genome, but which is highly variable between accessions, with one of our accessions carrying at least 6 copies (**Figure 4B).** However, there are no copies elsewhere in any of the new genomes (**Supplemental Figure 8).**

## Transposon-driven duplications

Transposons are thought to play a major role in gene duplications, moving genes or gene fragments around the genome (Woodhouse, Pedersen, and Freeling 2010). While confirming the *trans* duplications in the PacBio genomes, we found a beautiful example of this process. The gene AT1G20400 (annotated to encode a myosin heavy chain like protein, confirmed by expression data) was predicted to have multiple *trans*-duplications. The 944 bp coding region contains 125 putatively heterozygous SNPs with striking haplotype structure clearly suggesting structural variation **(Figure 4C)**. We were able to identify the duplication predicted by GWAS in the 6 PacBio genomes **(Figure 4).** Three of the newly assembled genomes have only one copy

180     of the gene, just like the reference genome, but one accession has 2 copies, and two have 4

181     copies. However, none of the 6 new genomes has a copy in the same place as in the reference

182     genome (**Supplemental Figure 9**).

183        In the reference genome, AT1G20400 is closely linked to AT1G203090, which is

184     annotated as a Gypsy transposable element. This element also contains many pseudo-SNPs,

185     and GWAS revealed duplication sites overlapping those for AT1G20400 (**Figure 3B**). This

186     suggested that the putative gene and putative Gypsy element transpose together, i.e. that both

187     are misannotated, and that the whole construct is effectively a large transposable element.

188     Further analysis of the PacBio genomes confirmed that AT1G20400 and AT1G20390 were

189     always found together, and we were also able to find conserved Long Terminal Repeat

190     sequences flanking the whole construct, as would be expected for a retrotransposon

191     (**Supplementary Figure 10-11**). Available bisulfite sequencing data (Kawakatsu et al. 2016)

192     showed that the whole region is heavily methylated, as expected for a transposon **(Figure 3).**

193     We tried mapping the bisulfite reads to the appropriate genome for the respective accesions,

194     but the coverage was too low and noisy to observe a difference in methylation between the

195     multiple insertions (**Supplemental Figure 12).**

**Figure 4:** A Gypsy element and a gene transpose together. (**A**) Methylation levels on regions containing AT1G20390 and AT1G20400 for 6 accessions, calculated in 200 bp windows after mapping reads to the TAIR10 reference genome (annotation shown in black). (**B**) GWAS results for the putatively heterozygous SNPs in AT1G20390 and AT1G20400. Each line represents the link between the position of the pseudo-SNP (upper side) and a GWAS hit position in the genome (lower side). The color corresponds to the location of the original heterozygous SNPs (phenotype), Blue for the gypsy element (AT1G20390) and orange for the myosin heavy chain-like gene (AT1G20400). The lower part shows the presence of the new transposable element in the 6 PacBio genomes. (**C**) Genotype around the AT1G20400 region for the 1001 genomes data set: Three different genotypes are shown, Homozygote reference (HOM REF), Heterozygot (HET) and Homozygote Alternative (HOM ALT). White represents a lack of coverage. I (**D**) Presence of the gene and the transposon in related species. Green represents the presence of the TE, gene or both in each species.

Having located precise insertions in the new genomes, we attempted to find them using short-read data in the 1001 Genomes dataset. Except for one insertion that was shared by 60% of accessions, the rest were found in less than 20% of accessions, suggesting that this new

211    element has no fixed insertions in the genome (**Supplemental Figure 13**). We looked for the

212    element of interest in the genomes of *A. lyrata* (two genomes), *A. suecica (Burns et al. 2021)*,

213    and the outgroup *Capsella rubella (Slotte et al. 2013)*. The gene and the Gypsy element were

214    only found together in *A. thaliana* (including the *A. thaliana* sub-genome of the allopolyploid *A.*

215    *suecica*). The Gypsy element alone is present in the other *Arabidopis* species, and the gene

216    alone is present in *A. lyrata*, but only in one of two genomes. In *Capsella rubella* neither the

217    transposon nor the gene could be detected **(Supplemental Figure 14)**. The transposon and

218    gene are specific to Arabidopsis while their co-transposition is specific to A. thaliana, suggesting

219    that these "jumps" occurred recently evolutionary speaking since divergence of A. thaliana with

220    the other Arabidopsis species.

# 221   Discussion

222    A duplication can lead to pseudo-SNPs when SNPs are identified by mapping short reads to a

223    reference genome that does not contain the duplication. Typically pseudo-SNPs have to be

224    identified using non-Mendelian segregation patterns in families or crosses, but in inbred lines

225    they can be identified solely by their presence. The overwhelming majority of the 3.3 million

226    heterozygous SNPs (44% of total) identified by our SNP-calling of the 1001 Genomes Project

227    (1001 Genomes Consortium 2016) data are likely to be pseudo-SNPs. Assuming this, we used

228    (pseudo-)heterozygosity as a "phenotype", and tried to map its cause, i.e. the duplication, using

229    a simple but powerful GWAS approach. Focusing on annotated genes, we find that over 2500

230    (roughly 10% of total) harbor pseudo-SNPs and show evidence of duplication. Using 6

231    independent long-read assemblies, we were able to confirm 60% of these duplications, using

232    conservative criteria (see Methods). Most of the remaining duplications are located in

233    pericentromeric regions where SNP-calling has lower quality, and which are difficult to assemble

234    even with long-read (**Supplemental Figure 15**).

235      These numbers nearly certainly underestimate the true extent of duplication. While

236    unlinked *trans*-duplications are fairly likely to give rise to pseudo-SNPs, local *cis*-duplications will

237    only do so once suffcent time has passed for substantial sequence divergence to occur. As for

238    the GWAS approach, it lacks statistical power to detect rare duplications, and can be misled by

239    allelic heterogeneity (due to multiple independent duplications). Finally, duplications are just a

240    subset of structural variants, and it is therefore not surprising that other short-read approaches

241    to detect such variants have identified many more (Zmienko et al. 2020; D.-X. Liu et al. 2021;

242    Göktay, Fulgione, and Hancock 2020).

243        Pseudo-SNPs is not the only problem with relying on a reference genome. Our analysis
244   uncovered a striking example of the potential importance of the "mobileome" in shaping genome
245   diversity: we show that an annotated gene and an annotated transposon are both part of a
246   much large mobile element, and the insertion in the reference genome is missing from most
247   other accessions. When short reads from another accession are mapped to this "gene" using
248   the reference genome, you are neither mapping to a gene, nor to the position you think.
249   Time (and more independently assembled genomes) will tell how significant this problem is, but
250   the potential for artifactual results is clearly substantial. It is also important to realize that the
251   artefactual nature of the 44% heterozygous SNPs was only apparent because we are working
252   with inbred lines. In human genetics, SNP-calling relies heavily on family trios, but in
253   outcrossing organisms where this is not possible, there is great cause for concern. The
254   increasing ease and ability to sequence more and more complex genomes, such as projects
255   associated with the 1001G+ and Tree of Life, will allow population analyses to avoid the use of
256   a single reference genome and reveal new mechanisms of gene duplication and structural
257   variants such as those reported here.

# 258   Methods

## 259   Genome assemblies

### 260   Long-read sequencing of six *A. thaliana*

261   We sequenced six Swedish *A. thaliana* lines part of the 1001 genomes collection ((1001
262   Genomes Consortium 2016)), ecotype ids 1254, 5856, 6021, 6024, 9412 and 9470. Plants
263   were grown in the growth chamber at 21 C in long day settings for 3 weeks and dark-treated for
264   24-48 hours before being collected. DNA was extracted from ~20 g of frozen whole seedling
265   material following a high molecular weight DNA extraction protocol adapted for plant tissue
266   (Cristina Barragan et al. 2021). All six genomes were sequenced with PacBio, 6021 with PacBio
267   RSII technology, while the rest were sequenced with Sequel. Accession 9412 had two rounds of
268   Sequel sequencing and 6024 was additionally sequenced with Nanopore (4.1 Gbp sequenced,
269   376 K reads with N50 18.7 Kb), all data was used in the assemblies.

## 270    MinION sequencing of two *A. lyrata*

271    We sequenced two North American *A. lyrata* accessions, 11B02 and 11B21. Both individuals
272    come from the 11B population of *A. lyrata*, which is self compatible and situated in Missouri
273    (Griffin and Willi 2014) (GPS coordinates 38° 28' 07.1" N; 90° 42' 34.3" W) . Plants were bulked
274    for 1 generation in the lab and DNA was extracted from ~20g of 3 week old seedlings, grown at
275    21°C and dark treated for 3 days prior to tissue collection. DNA was extracted using a modified
276    protocol for high molecular weight DNA extraction from plant tissue. DNA quality was assessed
277    with a Qubit fluorometer and a Nanodrop analysis. We used a Spot-ON Flow Cell FLO-
278    MIN106D R9 Version with a ligation sequencing kit SQK-LSK109. Bases were called using
279    guppy (https://nanoporetech.com/community) (version 3.2.6). The final output of MinION
280    sequencing for 11B02 was 13,67 Gbp in 763,800 reads and an N50 of 31,15 Kb. The final
281    output of MinION sequencing for 11B21 was 17.55 Gb, 1.11 M reads with an N50 of 33.26 Kb.

## 282    Genome assembly, polishing and scaffolding

283    The six *A. thaliana* genomes (ecotype ids 1254, 5856, 6021, 6024, 9412 and 9470) were
284    assembled using Canu (v 1.7.1) ((Koren et al. 2017)) with default settings, except for
285    genomeSize. Previous estimates of flow cytometry were used for this parameter (Long et al.
286    2013) when available or 170m was used. The values were 170m, 178m, 135m, 170m, 170m
287    and 170m, respectively. The assemblies were corrected with two rounds of arrow (PacBio's
288    SMRT Link software release 5.0.0.6792) and one of Pilon (Walker et al. 2014). For arrow, the
289    respective long reads were used and for Pilon the 1001 genomes DNA sequencing data, plus
290    PCR-free Illumina 150bp data that was generated for accessions 6024 and 9412; lines 5856,
291    6021, 9470 had available PCR-free data (250bp reads generated by David Jaffe, Broad
292    Institute). This resulted in 125.6Mb, 124.3Mb, 124.5Mb, 124.7Mb, 127.1Mb and 128Mb
293    assembled bases, respectively; contained in 99, 436, 178, 99, 109 and 124 contigs,
294    respectively. The polished contigs were ordered and scaffolded with respect to the Col-0
295    reference genome, using RaGOO (Alonge et al. 2019).
296    We assembled the genome of the two *A. lyrata* accessions 11B02 and 11B21 using
297    Canu (Koren et al. 2017) (v 1.8) with default settings and a genome size set to 200Mb. The
298    genome of 11B02 is contained in 498 contigs and the genome of 11B02 in 265 contigs. The
299    contig assemblies were polished using Racon (Vaser et al. 2017) (v 1.4) and ONT long reads
300    were mapped using nglmr (Sedlazeck et al. 2018) (v 0.2.7). Assemblies were further polished

300 by mapping PCR-free Illumina 150bp short reads (~100X for 11B02 and ~88X for 11B21) to the

301 long read corrected assemblies. Short read correction of assembly errors was carried out using

302 Pilon (Walker et al. 2014) (v1.23). Contigs were scaffolded into pseudo-chromosomes using

303 RaGOO (Alonge et al. 2019) and by using the error corrected long reads from Canu and the *A.*

304 *lyrata* reference genome (Hu et al. 2011)  and the *A. arenosa* subgenome of *A. suecica (Burns*

305 *et al. 2021)* as a guide followed by manual inspection of regions. The assembly size for 11B02

306 was 213Mb and 11B21 was 202Mb. Genome size was estimated using findGSE (Sun et al.

307 2018) with a resulting estimated genome size of ~256Mb for 11B02 and ~237Mb for 11B21.

## Heterozygous SNPs calling / extraction

309 We downloaded short-read data for 1,057 accessions from the 1001 Genomes Project (1001

310 Genomes Consortium 2016). Raw paired-end reads were processed with cutadapt (v1.9)

311 (Martin 2011) to remove 3' adapters, and to trim 5'-ends with quality 15 and 3'-ends with quality

312 10 or N-endings. All reads were aligned to the *A. thaliana* TAIR10 reference genome

313 (Arabidopsis Genome Initiative 2000) with BWA-MEM (v0.7.8) (H. Li 2013), and both Samtools

314 (v0.1.18) and Sambamba (v0.6.3) were used for various  file format conversions, sorting and

315 indexing (H. Li et al. 2009; Tarasov et al. 2015), while duplicated reads where by marked by

316 Markduplicates from Picard (v1.101; http://broadinstitute.github.io/picard/). Further steps were

317 carried out with GATK (v3.4) functions (Van der Auwera et al. 2013; DePristo et al. 2011). Local

318 realignment around indels were done with 'RealignerTargetCreator' and 'IndelRealigner', and

319 base recalibration with 'BaseRecalibrator' by providing known indels and SNPS from The 1001

320 Genomes Consortium (1001 Genomes Consortium 2016). Genetic variants were called with

321 'HaplotypeCaller' in individual samples followed by joint genotyping of a single cohort with

322 'GenotypeGVCFs'. An initial SNP filtering was done following the variant quality score

323 recalibration (VQSR) protocol. Briefly, a subset of ~181,000 high quality SNPs from the RegMap

324 panel (Horton et al. 2012) were used as the training set for VariantRecalibrator with a priori

325 probability of 15 and four maximum Gaussian distributions. Finally, only bi-allelic SNPs within at

326 a sensitivity tranche level of 99.5 were kept, for a total of 7,311,237 SNPs.

## SNPs filtering

From the raw VCF files only SNPs positions containing heterozygous labels have been extracted using GATK VariantFiltration. From the 3.3 millions of heterozygous SNPs extracted from the 1001 genome dataset, two filtering steps have been applied. First only SNPs with a frequency of at least 5% of the population have been kept. From those, only the one inside of annotated coding regions from the TAIR10 annotation have been used. After those filtering steps a core set of 26647 SNPs have been used for further analysis. **(see Supplemental Figure 2)** Genes name and features containing those Pseudo-SNPs have been extracted from the TAIR10 annotation.

## GWAS

The presence and absence of all pseudo-SNP from the core set (coded as 1 and 0 respectively) have been used as a phenotype to run GWAS. As a genotype the matrix published by the 1001 genome consortium containing 10 Millions SNPs has been used (1001 Genomes Consortium 2016). To run all the GWAS, the pygwas package (https://github.com/timeu/PyGWAS) with the amm (accelerated mixed model) option has been used. The raw output containing all SNPs has been filtered, removing all SNPs with a minor allele frequency below 0.05 and/ or a -log10(p-value) below 4.

For each GWAS performed, the p-value as well as the position have been used to call the peaks using the fourier transform function in R (filterFFT) combine with the peak detection function (peakDetection), from the package NucleR 3.13, to automatically retrieve the position of each peak across the genome. From each peak the Highest SNPs within a region of +/- 10kb around the peak center have been used. (Example presented in **Supplemental Figure 16)** Using all 26647 SNPs a summary table was generated with each pseudo-hete SNPs and each GWAS peak detected (**Supplemental Data**). This matrix was then used to generate **Figure 2C**, applying thresholds of -log10(pvalue) of 20 and minor allele frequency of 0.1.

## Confirmation of GWAS results

To confirm the different insertions detected a combination of blast and synteny was used on the denovo assembled genomes. Only the insertions that segregate in the 6 new genomes have been used (398). For each gene considered the corresponding sequence from the TAIR10

356  annotation got blast to each genomes. As genomes are not annotated, a blast of the genes
357  around the GWAS peak have been used to find the corresponding region in each genome, **see**
358  **Supplemental Figure 4**. The presence of a blast results within 20kb of the peak position have
359  been then used to confirm the GWAS results.

## Gene ontology

361  Out of the 2570 genes detected to be duplicated, 2396 have a gene ontology annotation.
362  PLAZA.4 (Van Bel et al. 2018) has been used to perform a gene enrichment analysis using the
363  full genome as background. Data has been then retrieved and plotted using R.

## Coverage and Methylation analysis

365  Bisulfite reads for the accessions were taken from 1001 methylomes (Kawakatsu et al. 2016).
366  Reads were mapped to PacBio genomes using an nf-core pipeline
367  (https://github.com/rbpisupati/methylseq). The weighted methylation levels are then later
368  calculated on windows of 200bp using custom python scripts (Schultz et al. 2012).
369  The sequencing coverage for each accession has been extracted using the function
370  bamCoverage (windows size of 50bp) from the program DeepTools (Ramírez et al. 2016). The
371  Bigwig files generated are then processed in R using the package rtracklayer. An overall
372  comparison between the mean sequencing coverage and the number of pseudo-SNPs detected
373  show that no correlation is observed. (**Supplemental Figure 18**)

## Multiple sequence alignment

375  For each insertion of the AT1G20390-AT1G20400 (Transposon-GENE) fragment, a fasta file
376  including 2kb on each side of the fragment was extracted, from each corresponding denovo
377  assembled genomes, using the getfasta function from bedtools (Quinlan and Hall 2010).
378  Multiple alignment was performed using KALIGN (Madeira et al. 2019). Visualisation and
379  comparison was done using Jalview 2 (Waterhouse et al. 2009).

## Structural variation analysis

381  To control the structure of the region around duplicated genes, the sequence from 3 genes up

382    and down from the gene of interest have been extracted. Each sequence has then been blasted

383    to each of the pacbio and the position of each blast results have been retrieved. The NCBI

384    BLAST (Altschul et al. 1990) have been used with a percentage of identity threshold of 70% and

385    all other parameters as default. From each blast results fragments with at least 50% of the input

386    sequence length have been selected and plotted using R.

387    ## Frequency of the insertions in the 1001 genomes dataset

388    The same sequences used for the multiple alignment have been used to confirm presence or

389    absence of each insertion in the 1001 genomes dataset. We used each of those sequences as

390    reference to map short reads using minimap 2 (H. Li 2018). For each insertion, only reads

391    having both pair mapping in the regions have been selected. An insertion has been validated as

392    present in an accession if at least 3 pairs of reads are spanning the insertion border. **(see**

393    **Supplemental Figure 10)**. For each insertion the frequency across the 1001 genome has been

394    extracted and presented in **Supplemental Figure 10.**

395    ## Multiple species comparison

396    Similarly as for the *A. thaliana* genomes, we used the *Capsella Rubella* and *A.arenosa*

397    published genomes (Slotte et al. 2013; Burns et al. 2021). In the case of *A. arenosa* we used

398    the subgenome of *A. suecica*. We blasted the TE-GENE fragments, extracted from the TAIR10

399    annotation, using the NCBI program as above. For *A.lyrata* two newly assembled genomes

400    were assembled using MinION sequencing. The presence of the transposon or the GENE

401    alone, the two together (full fragment) and the number of insertions have been extracted and

402    summed up in **Figure 4D**.

403    ## Singleton analysis:

404    From the 4.4 millions pseudo-hete SNPs, 1 millions singletons and 0.3 millions doubletons

405    pseudo-hete SNPs have been extracted. First, based on their positions we overlapped those

406    SNPs with the list of detected duplicated genes and compared with a thousand randomly

407    generated distribution of genes across the genome. We found that a highly significant 11%

408    overlap with genes detected to be duplicated.

409        For each Singleton segregating in the 6 Pacbio accessions the reads overlapping the

410   Pseudo-SNP have been extracted and re-mapped to the corresponding PacBio. The position,

411   the coverage as well and the number of SNPs detected have been extracted from the bam files.

# Acknowledgment

413   We thank Rahul Pisupati for providing the methylation data.

# Authors' contributions

415   BJ and MN developed the project. BJ performed all analyses. LMS and RB assembled the

416   *A.thaliana* and *A.lyrata* genomes, respectively. FR generated the SNP matrix. BJ and MN wrote

417   the manuscript, with input from all authors.

# Funding

419   This project has received funding from the European Research Council (ERC) under the

420   European Union's Horizon 2020 research and innovation programme (grant agreement No

421   789037)

# Availability of data and materials

423   All genome assemblies and raw reads were deposited under the BioProject ID: PRJNA779205.

424   Scripts used are available under Github link: https://github.com/benjj212/duplication-paper.git.

425   The full GWAS matrix is available at https://doi.org/10.5281/zenodo.5702395

# Ethics approval and consent to participate

427   Not applicable.

# 428 Competing interests

429 The authors declare no competing interests.

# 430 Author details

431 1 Gregor Mendel Institute, Austrian Academy of Sciences, Vienna Biocenter, Vienna, Austria. 2
432 Max Planck Institute for Developmental Biology, Tübingen, Germany. 3 Department of Plant
433 Sciences, University of Cambridge, Cambridge, UK.

# 434 References

435 1001 Genomes Consortium. 2016. "1,135 Genomes Reveal the Global Pattern of Polymorphism
436     in Arabidopsis Thaliana." *Cell* 166 (2): 481–91.

437 Alkan, Can, Bradley P. Coe, and Evan E. Eichler. 2011. "Genome Structural Variation Discovery
438     and Genotyping." *Nature Reviews. Genetics* 12 (5): 363–76.

439 Alonge, Michael, Sebastian Soyk, Srividya Ramakrishnan, Xingang Wang, Sara Goodwin, Fritz
440     J. Sedlazeck, Zachary B. Lippman, and Michael C. Schatz. 2019. "RaGOO: Fast and
441     Accurate Reference-Guided Scaffolding of Draft Genomes." *Genome Biology* 20 (1): 224.

442 Alonge, Michael, Xingang Wang, Matthias Benoit, Sebastian Soyk, Lara Pereira, Lei Zhang,
443     Hamsini Suresh, et al. 2020. "Major Impacts of Widespread Structural Variation on Gene
444     Expression and Crop Improvement in Tomato." *Cell*.
445     https://doi.org/10.1016/j.cell.2020.05.021.

446 Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. 1990. "Basic Local Alignment
447     Search Tool." *Journal of Molecular Biology* 215 (3): 403–10.

448 Arabidopsis Genome Initiative. 2000. "Analysis of the Genome Sequence of the Flowering Plant
449     Arabidopsis Thaliana." *Nature* 408 (6814): 796–815.

450 Burns, Robin, Terezie Mandáková, Joanna Gunis, Luz Mayela Soto-Jiménez, Chang Liu, Martin
451     A. Lysak, Polina Yu Novikova, and Magnus Nordborg. 2021. "Gradual Evolution of

452    Allopolyploidy in Arabidopsis Suecica." *Nature Ecology & Evolution* 5 (10): 1367–81.

453    Carter, Nigel P. 2007. "Methods and Strategies for Analyzing Copy Number Variation Using

454    DNA Microarrays." *Nature Genetics* 39 (7 Suppl): S16–21.

455    Cristina Barragan, A., Maximilian Collenberg, Rebecca Schwab, Merijn Kerstens, Ilja Bezrukov,

456    Felix Bemm, Doubravka Požárová, Filip Kolář, and Detlef Weigel. 2021. "Homozygosity at

457    Its Limit: Inbreeding Depression in Wild Arabidopsis Arenosa Populations." *bioRxiv*.

458    https://doi.org/10.1101/2021.01.24.427284.

459    DePristo, Mark A., Eric Banks, Ryan Poplin, Kiran V. Garimella, Jared R. Maguire, Christopher

460    Hartl, Anthony A. Philippakis, et al. 2011. "A Framework for Variation Discovery and

461    Genotyping Using next-Generation DNA Sequencing Data." *Nature Genetics* 43 (5): 491–

462    98.

463    Göktay, Mehmet, Andrea Fulgione, and Angela M. Hancock. 2020. "A New Catalogue of

464    Structural Variants in 1301 A. Thaliana Lines from Africa, Eurasia and North America

465    Reveals a Signature of Balancing at Defense Response Genes." *Molecular Biology and*

466    *Evolution*, November. https://doi.org/10.1093/molbev/msaa309.

467    Gonzalez, Enrique, Hemant Kulkarni, Hector Bolivar, Andrea Mangano, Racquel Sanchez,

468    Gabriel Catano, Robert J. Nibbs, et al. 2005. "The Influence of CCL3L1 Gene-Containing

469    Segmental Duplications on HIV-1/AIDS Susceptibility." *Science* 307 (5714): 1434–40.

470    Griffin, P. C., and Y. Willi. 2014. "Evolutionary Shifts to Self-Fertilisation Restricted to

471    Geographic Range Margins in North American Arabidopsis Lyrata." *Ecology Letters* 17 (4):

472    484–90.

473    Handsaker, Robert E., Joshua M. Korn, James Nemesh, and Steven A. McCarroll. 2011.

474    "Discovery and Genotyping of Genome Structural Polymorphism by Sequencing on a

475    Population Scale." *Nature Genetics* 43 (3): 269–76.

476    Horton, Matthew W., Angela M. Hancock, Yu S. Huang, Christopher Toomajian, Susanna

477    Atwell, Adam Auton, N. Wayan Muliyati, et al. 2012. "Genome-Wide Patterns of Genetic

478    Variation in Worldwide Arabidopsis Thaliana Accessions from the RegMap Panel." *Nature*

479    *Genetics* 44 (2): 212–16.

480    Hufford, Matthew B., Arun S. Seetharam, Margaret R. Woodhouse, Kapeel M. Chougule,

481    Shujun Ou, Jianing Liu, William A. Ricci, et al. 2021. "De Novo Assembly, Annotation, and

482    Comparative Analysis of 26 Diverse Maize Genomes." *Cold Spring Harbor Laboratory*.

483    https://doi.org/10.1101/2021.01.14.426684.

484    Hu, Tina T., Pedro Pattyn, Erica G. Bakker, Jun Cao, Jan-Fang Cheng, Richard M. Clark, Noah

485    Fahlgren, et al. 2011. "The Arabidopsis Lyrata Genome Sequence and the Basis of Rapid

486     Genome Size Change." *Nature Genetics* 43 (5): 476–81.

487  Jiao, Wen-Biao, and Korbinian Schneeberger. 2019. "Chromosome-Level Assemblies of

488     Multiple Arabidopsis Thaliana Accessions Reveal Hotspots of Genomic Rearrangements."

489     *bioRxiv*. https://doi.org/10.1101/738880.

490  Kawakatsu, Taiji, Shao-Shan Carol Huang, Florian Jupe, Eriko Sasaki, Robert J. Schmitz, Mark

491     A. Urich, Rosa Castanon, et al. 2016. "Epigenomic Diversity in a Global Collection of

492     Arabidopsis Thaliana Accessions." *Cell* 166 (2): 492–505.

493  Koren, Sergey, Brian P. Walenz, Konstantin Berlin, Jason R. Miller, Nicholas H. Bergman, and

494     Adam M. Phillippy. 2017. "Canu: Scalable and Accurate Long-Read Assembly via Adaptive

495     K-Mer Weighting and Repeat Separation." *Genome Research* 27 (5): 722–36.

496  Li, Changsheng, Xiaoli Xiang, Yongcai Huang, Yong Zhou, Dong An, Jiaqiang Dong, Chenxi

497     Zhao, et al. 2020. "Long-Read Sequencing Reveals Genomic Structural Variations That

498     Underlie Creation of Quality Protein Maize." *Nature Communications* 11 (1): 17.

499  Li, Heng. 2013. "Aligning Sequence Reads, Clone Sequences and Assembly Contigs with BWA-

500     MEM." *arXiv [q-bio.GN]*. arXiv. http://arxiv.org/abs/1303.3997.

501  ———. 2018. "Minimap2: Pairwise Alignment for Nucleotide Sequences." *Bioinformatics*  34

502     (18): 3094–3100.

503  Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth,

504     Goncalo Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup.

505     2009. "The Sequence Alignment/Map Format and SAMtools." *Bioinformatics*  25 (16):

506     2078–79.

507  Lin, Ke, Ningwen Zhang, Edouard I. Severing, Harm Nijveen, Feng Cheng, Richard G. F.

508     Visser, Xiaowu Wang, Dick de Ridder, and Guusje Bonnema. 2014. "Beyond Genomic

509     Variation - Comparison and Functional Annotation of Three Brassica Rapagenomes: A

510     Turnip, a Rapid Cycling and a Chinese Cabbage." *BMC Genomics* 15 (1): 250.

511  Liu, Dong-Xu, Ramesh Rajaby, Lu-Lu Wei, Lei Zhang, Zhi-Quan Yang, Qing-Yong Yang, and

512     Wing-Kin Sung. 2021. "Calling Large Indels in 1047 Arabidopsis with IndelEnsembler."

513     *Nucleic Acids Research*, October. https://doi.org/10.1093/nar/gkab904.

514  Liu, Yucheng, Huilong Du, Pengcheng Li, Yanting Shen, Hua Peng, Shulin Liu, Guo-An Zhou, et

515     al. 2020. "Pan-Genome of Wild and Cultivated Soybeans." *Cell*, June.

516     https://doi.org/10.1016/j.cell.2020.05.023.

517  Long, Quan, Fernando A. Rabanal, Dazhe Meng, Christian D. Huber, Ashley Farlow, Alexander

518     Platzer, Qingrun Zhang, et al. 2013. "Massive Genomic Variation and Strong Selection in

519     Arabidopsis Thaliana Lines from Sweden." *Nature Genetics* 45 (8): 884–90.

520   Madeira, Fábio, Young Mi Park, Joon Lee, Nicola Buso, Tamer Gur, Nandana Madhusoodanan,
521       Prasad Basutkar, et al. 2019. "The EMBL-EBI Search and Sequence Analysis Tools APIs in
522       2019." *Nucleic Acids Research* 47 (W1): W636–41.
523   Martin, Marcel. 2011. "Cutadapt Removes Adapter Sequences from High-Throughput
524       Sequencing Reads." *EMBnet.journal* 17 (1): 10–12.
525   Miyahara, E., J. Pokorny, V. C. Smith, R. Baron, and E. Baron. 1998. "Color Vision in Two
526       Observers with Highly Biased LWS/MWS Cone Ratios." *Vision Research* 38 (4): 601–12.
527   Perry, George H., Nathaniel J. Dominy, Katrina G. Claw, Arthur S. Lee, Heike Fiegler, Richard
528       Redon, John Werner, et al. 2007. "Diet and the Evolution of Human Amylase Gene Copy
529       Number Variation." *Nature Genetics* 39 (10): 1256–60.
530   Pinosio, Sara, Stefania Giacomello, Patricia Faivre-Rampant, Gail Taylor, Veronique Jorge,
531       Marie Christine Le Paslier, Giusi Zaina, et al. 2016. "Characterization of the Poplar Pan-
532       Genome by Genome-Wide Identification of Structural Variation." *Molecular Biology and*
533       *Evolution* 33 (10): 2706–19.
534   Quadrana, Leandro, Amanda Bortolini Silveira, George F. Mayhew, Chantal LeBlanc, Robert A.
535       Martienssen, Jeffrey A. Jeddeloh, Vincent Colot, and Daniel Zilberman. 2016. "The
536       Arabidopsis Thaliana Mobilome and Its Impact at the Species Level." *eLife* 5 (June):
537       e15716.
538   Quinlan, Aaron R., and Ira M. Hall. 2010. "BEDTools: A Flexible Suite of Utilities for Comparing
539       Genomic Features." *Bioinformatics*  26 (6): 841–42.
540   Ramírez, Fidel, Devon P. Ryan, Björn Grüning, Vivek Bhardwaj, Fabian Kilpert, Andreas S.
541       Richter, Steffen Heyne, Friederike Dündar, and Thomas Manke. 2016. "deepTools2: A next
542       Generation Web Server for Deep-Sequencing Data Analysis." *Nucleic Acids Research* 44
543       (W1): W160–65.
544   Ranade, K., M. S. Chang, C. T. Ting, D. Pei, C. F. Hsiao, M. Olivier, R. Pesich, et al. 2001.
545       "High-Throughput Genotyping with Single Nucleotide Polymorphisms." *Genome Research*
546       11 (7): 1262–68.
547   Sedlazeck, Fritz J., Philipp Rescheneder, Moritz Smolka, Han Fang, Maria Nattestad, Arndt von
548       Haeseler, and Michael C. Schatz. 2018. "Accurate Detection of Complex Structural
549       Variations Using Single-Molecule Sequencing." *Nature Methods* 15 (6): 461–68.
550   Shendure, Jay, and Hanlee Ji. 2008. "Next-Generation DNA Sequencing." *Nature Biotechnology*
551       26 (10): 1135–45.
552   Slotte, Tanja, Khaled M. Hazzouri, J. Arvid Ågren, Daniel Koenig, Florian Maumus, Ya-Long
553       Guo, Kim Steige, et al. 2013. "The Capsella Rubella Genome and the Genomic

554     Consequences of Rapid Mating System Evolution." *Nature Genetics* 45 (7): 831–35.

555  Snijders, A. M., N. Nowak, R. Segraves, S. Blackwood, N. Brown, J. Conroy, G. Hamilton, et al.

556     2001. "Assembly of Microarrays for Genome-Wide Measurement of DNA Copy Number."

557     *Nature Genetics* 29 (3): 263–64.

558  Sun, Hequan, Jia Ding, Mathieu Piednoël, and Korbinian Schneeberger. 2018. "findGSE:

559     Estimating Genome Size Variation within Human and Arabidopsis Using K-Mer

560     Frequencies." *Bioinformatics*  34 (4): 550–57.

561  Tarasov, Artem, Albert J. Vilella, Edwin Cuppen, Isaac J. Nijman, and Pjotr Prins. 2015.

562     "Sambamba: Fast Processing of NGS Alignment Formats." *Bioinformatics*  31 (12): 2032–

563     34.

564  Van Bel, Michiel, Tim Diels, Emmelien Vancaester, Lukasz Kreft, Alexander Botzki, Yves Van

565     de Peer, Frederik Coppens, and Klaas Vandepoele. 2018. "PLAZA 4.0: An Integrative

566     Resource for Functional, Evolutionary and Comparative Plant Genomics." *Nucleic Acids*

567     *Research* 46 (D1): D1190–96.

568  Van der Auwera, Geraldine A., Mauricio O. Carneiro, Chris Hartl, Ryan Poplin, Guillermo Del

569     Angel, Ami Levy-Moonshine, Tadeusz Jordan, et al. 2013. "From FastQ Data to High

570     Confidence Variant Calls: The Genome Analysis Toolkit Best Practices Pipeline." *Current*

571     *Protocols in Bioinformatics / Editoral Board, Andreas D. Baxevanis ... [et Al.]* 11 (1110):

572     11.10.1–11.10.33.

573  Vaser, Robert, Ivan Sović, Niranjan Nagarajan, and Mile Šikić. 2017. "Fast and Accurate de

574     Novo Genome Assembly from Long Uncorrected Reads." *Genome Research* 27 (5): 737–

575     46.

576  Walker, Bruce J., Thomas Abeel, Terrance Shea, Margaret Priest, Amr Abouelliel, Sharadha

577     Sakthikumar, Christina A. Cuomo, et al. 2014. "Pilon: An Integrated Tool for

578     Comprehensive Microbial Variant Detection and Genome Assembly Improvement." *PloS*

579     *One* 9 (11): e112963.

580  Waterhouse, Andrew M., James B. Procter, David M. A. Martin, Michèle Clamp, and Geoffrey J.

581     Barton. 2009. "Jalview Version 2--a Multiple Sequence Alignment Editor and Analysis

582     Workbench." *Bioinformatics*  25 (9): 1189–91.

583  Woodhouse, Margaret R., Brent Pedersen, and Michael Freeling. 2010. "Transposed Genes in

584     Arabidopsis Are Often Associated with Flanking Repeats." *PLoS Genetics* 6 (5): e1000949.

585  Yao, Wen, Guangwei Li, Hu Zhao, Gongwei Wang, Xingming Lian, and Weibo Xie. 2015.

586     "Exploring the Rice Dispensable Genome Using a Metagenome-like Assembly Strategy."

587     *Genome Biology* 16 (September): 187.

588    Zhao, Min, Qingguo Wang, Quan Wang, Peilin Jia, and Zhongming Zhao. 2013. "Computational

589        Tools for Copy Number Variation (CNV) Detection Using next-Generation Sequencing

590        Data: Features and Perspectives." *BMC Bioinformatics* 14 Suppl 11 (September): S1.

591    Zhou, Yong, Dmytro Chebotarov, Dave Kudrna, Victor Llaca, Seunghee Lee, Shanmugam

592        Rajasekar, Nahed Mohammed, et al. 2020. "A Platinum Standard Pan-Genome Resource

593        That Represents the Population Structure of Asian Rice." *Scientific Data* 7 (1): 113.

594    Zmienko, Agnieszka, Malgorzata Marszalek-Zenczak, Pawel Wojciechowski, Anna Samelak-

595        Czajka, Magdalena Luczak, Piotr Kozlowski, Wojciech M. Karlowski, and Marek

596        Figlerowicz. 2020. "AthCNV: A Map of DNA Copy Number Variations in the Arabidopsis

597        Genome." *The Plant Cell* 32 (6): 1797–1819.