1  **The nature of intraspecific genome size variation in taxonomically complex eyebrights**

2  **Authors:**

3  Hannes Becher[1], Robyn F. Powell[2], Max R. Brown[1,3], Chris Metherell[4], Jaume Pellicer[2,5], Ilia J.
4  Leitch[2], Alex D. Twyford[1,6,7]

5  1 Institute of Evolutionary Biology, School of Biological Sciences, University of Edinburgh,
6  Charlotte Auerbach Road, Edinburgh, EH9 3FL, UK

7  2 Jodrell Laboratory, Royal Botanic Gardens, Kew, Richmond, Surrey, UK

8  3 current address: Wellcome Trust Genome Campus, Hinxton, Saffron Walden CB10 1RQ

9  4 Botanical Society of Britain and Ireland, 4 High Firs Crescent, Harpenden, Hertfordshire AL5
10  1NA, UK

11  5 Institut Botànic de Barcelona (IBB, CSIC-Ajuntament de Barcelona), Passeig del Migdia sn,
12  08038, Barcelona, Spain

13  6 Royal Botanic Garden Edinburgh, 20A Inverleith Row, Edinburgh EH3 5LR, UK

14  7 correspondence: Alex.Twyford@ed.ac.uk

15  Word counts:

16  Total word count for the main body of the text (Introduction, Materials and Methods, Results,
17  Discussion):  5826

18  Section word counts: Into 1035, M&M 1320, Results 1730, Discussion 1741

19  Figures: 3, all colour

20  Tables: 1

21  Supporting Figures: 2

22  Supporting Tables: 2

23

24

25

26

**Summary**

- Genome size (GS) is a key trait related to morphology, life history, and evolvability. Although GS is, by definition, affected by presence/absence variants (PAVs), which are ubiquitous in population sequencing studies, GS is often treated as an intrinsic property of a species. Here, we studied intra- and interspecific GS variation in taxonomically complex British eyebrights (*Euphrasia*).

- We generated GS data for 192 individuals of diploid and tetraploid *Euphrasia* and analysed GS variation in relation to ploidy, taxonomy, population affiliation, and geography. We further compared the genomic repeat content of 30 samples.

- We found considerable genuine intraspecific GS variation, and observed isolation-by-distance for GS in outcrossing diploids. Tetraploid *Euphrasia* showed contrasting patterns, with GS increasing with latitude in outcrossing *Euphrasia arctica*, but little GS variation in the highly selfing *Euphrasia micrantha*. Interspecific differences in GS genomic repeat percentages were small.

- We show the utility of treating GS as the outcome of polygenic variation. Like other types of genetic variation, such as single nucleotide polymorphisms, GS variation may be increased through hybridisation and population subdivision. In addition to selection on associated traits, GS is predicted to be affected indirectly by selection due to pleiotropy of the underlying PAVs.

**Keywords:** Genome size, polygenic trait, *Euphrasia*, ploidy, intraspecific variation, selection, pleiotropy, genomic repeats

49 **INTRODUCTION**

50 Genome size (GS), defined as the amount of DNA in an individual's unreplicated haploid

51 nucleus (Greilhuber *et al.*, 2005), is associated with an organisms life history strategy,

52 development, physiology, ecology, and gene and genome dynamics and evolution (Van't Hof &

53 Sparrow, 1963; Beaulieu *et al.*, 2008; Šímová & Herben, 2012; Greilhuber & Leitch, 2013;

54 Simonin & Roddy, 2018; Bilinski *et al.*, 2018; Novák *et al.*, 2020a; Roddy *et al.*, 2020). Genome

55 size is estimated to show a c. 64,000-fold variation across Eukaryotes, and c. 2440-fold

56 variation in flowering plants (Pellicer *et al.*, 2018). Much is known about broad-scale variation in

57 GS across land plants and algae, with different phyla characterised by different GS ranges

58 (Pellicer & Leitch, 2020), and showing, in many cases, a strong phylogenetic signal (e.g. Weiss-

59 Schneeweiss *et al.*, 2006; Vallès *et al.*, 2013; Wang *et al.*, 2016; Bainard *et al.*, 2019; Cacho *et

60 al.*, 2021). Genome size has also been explored in polyploid species, with studies showing that

61 while whole genome duplication events initially lead to an increase in GS, their subsequent

62 evolution is often accompanied by genome downsizing over time (Leitch *et al.*, 2008; Leitch &

63 Leitch, 2008; Pellicer *et al.*, 2010; Wong & Murray, 2012; Wendel, 2015; Zenil-Ferguson *et al.*,

64 2016). Recently, community ecology studies have started to include data on GS and

65 demonstrate its influence in shaping plant diversity (Guignard *et al.*, 2016, 2019). While

66 representative GS estimates have been obtained for approximately two thirds of flowering plant

67 families (Pellicer & Leitch, 2020), variation between individuals and populations within species

68 has typically received less attention, despite the increasing realisation that such variation may

69 be common (e.g. Šmarda *et al.*, 2010; Kolář *et al.*, 2017).

70 GS has often been considered a property of a species, and there has been much debate as to

71 whether it genuinely varies within species (Greilhuber, 2005; Gregory & Johnston, 2008;

72 Šmarda & Bureš, 2010). Genuine intraspecific differences in DNA content have been reported

73 or are predicted between individuals with: (1) heteromorphic sex chromosomes (Costich *et al.*,

74 1991; Renner *et al.*, 2017), (2) different numbers of B chromosomes (Leitch *et al.*, 2007)

75 dysploidy and aneuploidy, or (3) the presence/absence of specific DNA sequences such as (a)

76 structural variants including insertion-deletion polymorphisms (indels), (b) copy number variation

77 in protein-coding genes, commonly found in pan-genome studies (Hirsch *et al.*, 2014; Wang *et

78 al.*, 2018b; Gao *et al.*, 2019; Hübner *et al.*, 2019; Göktay *et al.*, 2020), and (c) copy number

79 variation of rDNA copies (Long *et al.*, 2013) or of other genomic repeats (Chia *et al.*, 2012;

80 Haberer *et al.*, 2020). Some differences, such as small indels, can be as small as one base pair,

81 while others are large-scale (many megabases), including sequence duplications or loss of a

82    dispensable chromosome. The above types of genetic variation can be subsumed under the

83    term presence/absence variants (PAVs), a type of structural genomic variation, and may be

84    detectable by methods for estimating GS, such as flow cytometry. Modern protocols using flow

85    cytometry with appropriate reference standards, and following best practice approaches, can be

86    accurate and highly precise (Greilhuber *et al.*, 2007; Pellicer *et al.*, 2021) and reveal genuine

87    intraspecific variation that can be confirmed by genome sequencing. Such sequencing has also

88    been used to reveal that repeat differences can be useful genetic markers, including

89    microsatellites and AFLPs. Consequently, there are an increasing number of well-documented

90    reports of genuine intraspecific GS variation (e.g. Achigan-Dako *et al.*, 2008; Šmarda *et al.*,

91    2010; Díez *et al.*, 2013; Hanušová *et al.*, 2014; Blommaert, 2020).


92    Our study considers such variation as polygenic, meaning heritable, and with a value affected

93    by multiple independent loci in the genome (Figure **1**). Many polygenic traits are known, with the

94    most renowned example being human height (Fisher, 1919). Loci underpinning polygenic

95    variation need not be protein-coding genes, but may also involve non-coding sequences

96    including introns, promotors, trans elements, or genomic repeats. Loci underpinning a polygenic

97    trait may differ in their effect sizes, as shown by Koornneef *et al.* (1991) for flowering time in

98    *Arabidopsis thaliana* (see also Napp-Zinn, 1955). Further, variants at a genetic locus are

99    commonly pleiotropic, affecting multiple traits and thus potentially being the target of multiple

100   selective effects. An early example of treating GS as such is the study of Meagher *et al.* (2005)

101   on the relationship between GS and flower size in *Silene latifolia*, that showed correlations

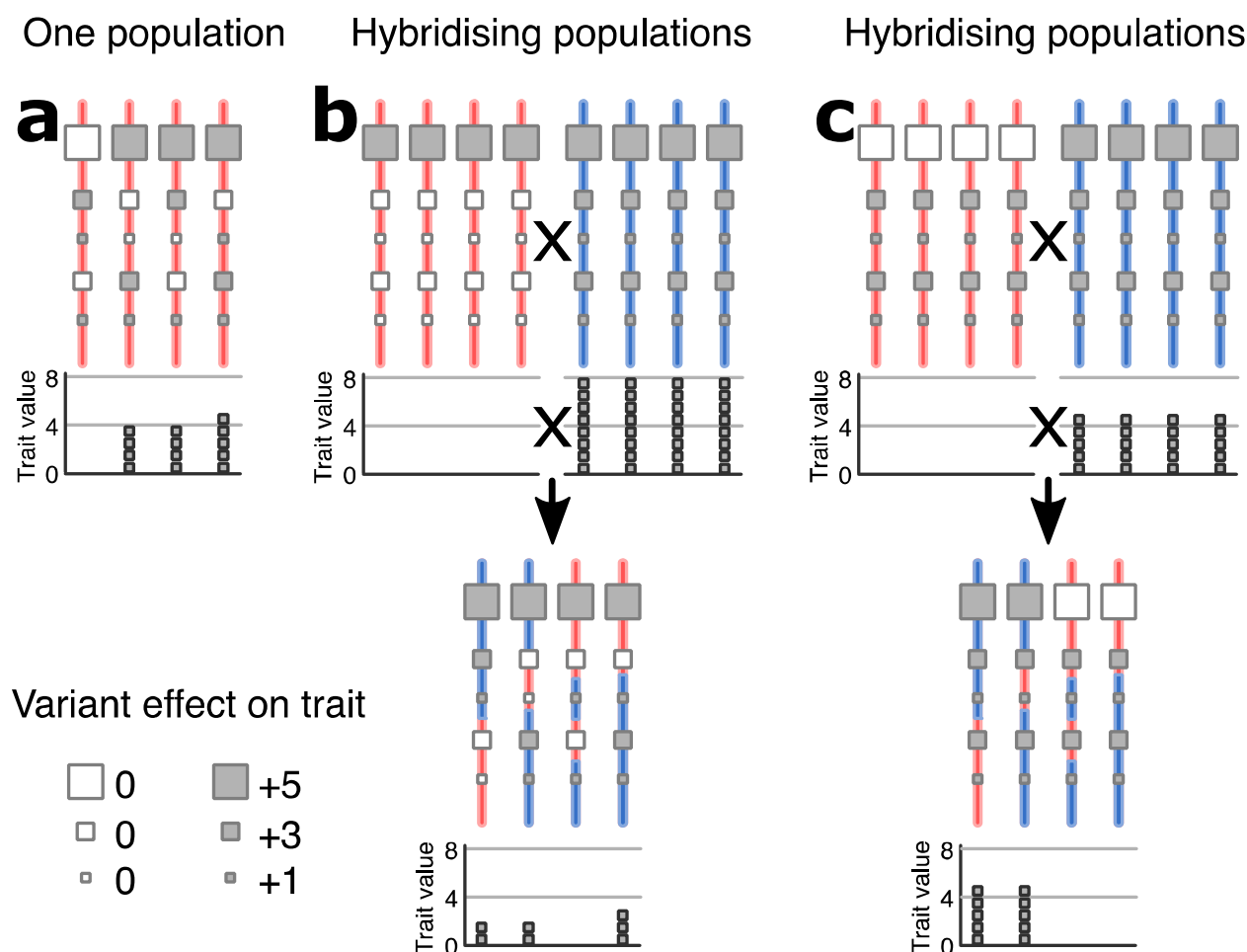102   between floral traits and GS in male plants.

**Figure 1. Schematic illustration of a polygenic trait, and its variability after hybridisation.**
Each red or blue line represents an individual's genome. Squares represent genetic variants with different effect sizes on a trait. The bar charts indicate individuals' trait values, relative to the individual with the lowest value. (a) A population (or species) with genetic variability for the trait. The effect of hybridisation between populations with different trait values depends on the genetic architecture of the trait difference. If the populations differ in many variants with small effects (b), recombinant offspring (denoted by mixed red and blue lines) are likely to have similar trait values. If, however, trait differences are due to a few variants with large effects (c), segregation in the recombinant offspring can produce higher trait variation. Applied to GS, open squares correspond to DNA missing and filled squares to DNA present at some site in the genome, as detailed in the main text.

5

116  Here we explore GS variation in British eyebrights (*Euphrasia* L., Orobanchaceae), a recently

117  radiating group of facultative hemiparasites. They comprise five diploid (2n = 2x = 22) and 15

118  tetraploid species (2n = 4x = 44) (Metherell & Rumsey, 2018), with genomic sequencing

119  showing that British tetraploids are closely related allotetraploids, with one sub-genome derived

120  from, or closely related to, British diploids (Becher *et al.*, 2020). The genus is an ideal model

121  group for investigating GS variation within and between closely related species because species

122  diversification is frequently postglacial (Gussarova *et al.*, 2008; Wang *et al.*, 2018a), with many

123  taxa being narrow endemics or recent hybrid species. *Euphrasia* therefore provides multiple

124  opportunities to study GS changes at the early stages of species divergence. Heterogeneous

125  ecological conditions may promote local adaptation, and extensive hybridisation may result in

126  local geographic homogenisation with variation in GS structured by geography rather than by

127  taxonomy, as seen previously in microsatellite and AFLP studies of population structure

128  (Kolseth & Lönn, 2005; French *et al.*, 2008).

129  To investigate the nature of GS variation in British *Euphrasia* species, we generated a

130  comprehensive dataset of 192 GS estimates across 13 species and 10 hybrid combinations,

131  supplemented with genomic sequence data to estimate the abundance of genomic repeats for

132  30 diverse diploids and tetraploids. Our study aims to answer the following questions: (1) How

133  variable is GS within species, between species, and between ploidy levels? (2) What is the

134  contribution of genomic repeats to GS variation in British *Euphrasia*, and how does repeat

135  content differ between the ploidy levels? (3) Does GS variation correspond with known patterns

136  of genetic structure and/or environmental variables in British *Euphrasia*? We discuss our results

137  in the light of polygenic variation, and we argue for a closer integration of population genomics

138  research with research on GS variation.

139

140  **METHODS**

141  **Population and species-level genome size variation**

142  **Population sampling.** Our sampling for GS estimation aimed to collect from across the

143  diversity of British *Euphrasia* taxa, and from a wide geographic area. Samples from 90

144  populations comprising 13 species and 10 hybrid combinations were either wild-collected and

145  used directly for GS estimates (54 samples) or collected as seeds and grown at the Royal

146  Botanic Garden Edinburgh following Brown *et al.* (2020) prior to GS estimation (138 samples). A

6

147　full list of samples analysed including their origin is given in Supplementary Information Table

148　S1. The identification of species and hybrids were made by the *Euphrasia* taxonomic expert

149　Chris Metherell, based on morphology.

150　**Genome size measurements.** Nuclear DNA content of *Euphrasia* samples was estimated by

151　flow cytometry using propidium iodide (PI) stained nuclei, following the one step method (see

152　Pellicer *et al.*, 2021). Briefly, for each *Euphrasia* accession, two small leaves (c. 1-2 cm) were

153　chopped together with the internal standard *Oryza sativa* 'IR36' (1C = 0.5 pg; Bennett & Smith,

154　1991) using a new razor blade, in a petri dish containing 1 mL of 'general purpose isolation

155　buffer' (GPB; Loureiro *et al.*, 2007), supplemented with 3% PVP-40 and 0.4 µL of β-

156　mercaptoethanol. An additional 1 mL of buffer was added to the homogenate, and then this was

157　filtered through a 30 µm nylon mesh to discard debris. Finally, the sample was stained with

158　100 µl of PI (1 mg/mL, Sigma) and incubated for 20 min on ice. For each accession analysed,

159　one sample was prepared, and this was run three times on the flow cytometer. The nuclear

160　DNA content of each sample run was estimated by recording at least 5,000 particles (c.1,000

161　nuclei per fluorescence peak) using a Cyflow SL3 flow cytometer (Sysmex-Partec GmbH,

162　Munster, Germany) fitted with a 100-mW green solid-state laser (Cobolt Samba). Resulting

163　output histograms were analysed using the FlowMax software (v. 2.9, Sysmex-Partec GmbH)

164　for statistical calculations. We report only GS estimates for samples where the coefficients of

165　variation (CV) of the sample and standard peaks in the flow histogram were less than 5% (see

166　Supporting Information Figure S1a and b for illustrative histograms of each ploidy level).

167　Where differences in GS were detected within a species, combined samples containing at least

168　two accessions were prepared following the same procedure as for individual runs. Genuine

169　intraspecific variation was confirmed where multiple fluorescence peaks were identified from the

170　combined run.

171　Throughout the paper we give 1C values in pg, where necessary converting published GS

172　values reported in Gbp to pg using a conversion factor of 0.978 following Doležel *et al.* (2003).

173　**Repeat content variation**

174　**Sequence data generation.** We used a combination of existing and newly generated genomic

175　sequencing data to investigate repeat variation in 31 samples comprising seven diploids and 23

176　tetraploids of *Euphrasia* plus *Bartsia alpina* as an outgroup. We downloaded short-read Illumina

177　data from the sequence read archive (SRA, see Supplementary Information Table S2). These

178    included 18 samples in total, including 12 tetraploid samples from the isolated island of Fair Isle

179    (Shetland, Scotland) generated for the study of Becher *et al.* (2020), which allowed us to study

180    genomic repeat profiles in sympatric populations. This dataset also included a total of six

181    representative diploid and tetraploid species from elsewhere in Britain.

182    We supplemented this previous data with newly generated sequence data from eleven

183    additional UK samples representing a wider range of species and geographic locations,

184    including 11 UK *Euphrasia* samples, an Austrian sample of *Euphrasia cuspidata* intended as a

185    close outgroup to UK species, and *Bartsia alpina* as an outgroup to the full sample set. Genomic

186    DNA was extracted from 12 silica-dried samples and herbarium material of *E. cuspidata* using

187    the Qiagen Plant Mini Kit (Qiagen, Manchester, UK), and used to prepare NEBUltra PCR-based

188    libraries. Pooled libraries were sent to Edinburgh Genomics where they were run on a single

189    lane of HiSeq 2500 using high output mode with 125 bp paired-end sequencing.

190    **Repeat content.** We ran the RepeatExplorer2 (RE) pipeline (https://repeatexplorer-elixir.cerit-

191    sc.cz/; Novák *et al.*, 2010, 2013, 2020) on a data set of 25,000 randomly selected read pairs of

192    each of the 31 samples (1,550,000 reads in total). This slightly exceeded the maximal number

193    of reads that can be analysed with default settings (which depends on the data). Our dataset

194    was therefore down-sampled to approximately 20,500 read pairs per sample. In comparative RE

195    analyses, read numbers are often supplied in proportion to genome sizes to assure repeats of

196    similar genome proportion can be detected in all samples (e.g. Novák *et al.*, 2020a). This logic

197    does not apply here, where the British samples comprise 23 closely related tetraploids and six

198    closely related diploids, with the diploid genome very similar to one of the tetraploid sub-

199    genomes (Becher *et al.*, 2020). No matter what genome proportion is chosen per sample, there

200    will always be more of the shared sub-genome than of the sub-genome restricted to tetraploids.

201    To minimise mate overlaps of short insert sizes, each read was trimmed to 100 nucleotides.

202    Further, we only used reads where at least 90 nucleotides had phred quality scores > 30. To

203    analyse the genomic repeat content, we excluded clusters annotated by RE as plastid DNA or

204    Illumina process controls. Our numbers thus deviate slightly from RE's automatic annotation.

205    **Statistical analyses.** Most GS analyses were conducted across all individuals or populations.

206    However, for *E. arctica*, *E. anglica*, and *E. micrantha*, where sampling covered most of their

207    large geographical range in Britain, we also analysed data from each species separately. All

208    analyses were done using R version 3.6.1 (R Core Team, 2019). For analyses of variance

209    (ANOVAs) we used the function aov(). To test whether sample means of GS were significantly

8

210   different, we used the function t.test(), with Bonferroni correction in cases of multiple testing. To

211   analyse how GS variation was partitioned by ploidy, taxon, and population we used ANOVA. To

212   test the effect of 'species', we then re-ran the ANOVAs without hybrids (Table **1**). To test the

213   significance of GS variance differences between species pairs, we divided the population mean

214   genome sizes by each species' grand mean (centring) and then applied an *F* test (R function

215   var.test()).

216   We tested the association between GS and latitude using a mixed effect model (R package

217   nlme, function lme()). For species analysed separately, we used linear models. We carried out

218   Mantel tests to assess the relationship between geographic distance and GS difference across

219   all samples as done by Duchoslav *et al.* (2013). Unlike genetic data, which require population

220   information, these Mantel tests could be carried out on individual-based genome size

221   differences or population means. Isolation by distance was assessed using Mantel tests (R

222   package vegan version 2.5-6) with 999 permutations.

223   To analyse genomic repeat patterns, we used hierarchical clustering and PCA on a matrix of the

224   per-sample genome proportions of the 100 largest repeat clusters in R using the functions

225   hclust() and prcomp(). *Bartsia alpina* was removed from the final PCA data set, because its

226   divergence from *Euphrasia* accounted for most of the variance in the data set, obscuring

227   variation within *Euphrasia*. To identify repeat clusters with large contributions to the first

228   principal component, we selected those clusters which had absolute values > 0.1 in the first

229   eigenvector. We further used binomial-family generalised linear models to estimate the average

230   genomic proportion individually for each repeat cluster. For each estimate, we computed the

231   residual sum of squares as a measure of the variation in genomic abundance between

232   individuals. We used linear models to assess the differences in relative abundance of individual

233   repeat types between ploidy levels.

234   To investigate a possible association of individual repeat clusters with GS, we used nine

235   tetraploid samples for which we had both an estimate of the population average GS and repeat

236   data (samples marked with asterisks in Supporting Information Table S2). We used the function

237   cor.test() to assess the significance level of any associations between the genome proportion of

238   each individual repeat cluster and population average GS.

239

240   **Results**

9

241    *Population and species-level genome size variation*

242    Genome size estimates from all 192 individuals passed our quality checks. These samples

243    came from 13 different species and 10 hybrid combinations, including 40 diploid and 152

244    tetraploid individuals (Supporting Information Table **S1**). Our samples covered a particularly

245    wide geographic range for the large-flowered species *E. anglica* (diploid, 552 km) and *E. arctica*

246    (tetraploid, 1152 km), and the small-flowered and highly selfing *E. micrantha* (tetraploid,

247    962 km).

248    The mean GS across all tetraploids was 1.18 pg (s.e. 0.004 pg), which is 11% less than twice

249    the mean GS of the diploids (0.66 pg, s.e. 0.008 pg). In the diploids, individual values ranged

250    1.2-fold, from 0.60 pg in *E. anglica* (population BED) to 0.73 pg in *E. anglica* in Dumfriesshire

251    (E4E0085). In tetraploids there was 1.3-fold variation, from 0.99 pg in *E. foulaensis* in Fair Isle

252    (FIA105) to 1.33 pg in *E. arctica* in Orkney (E4E0033).

253    Intraspecific GS ranges were widest in *E. arctica* (n = 43) and *E. foulaensis* (n = 13) (both 1.3-

254    fold), and *E. anglica* (n = 23) (1.2-fold). *Euphrasia confusa* (n = 6), *E. nemorosa* (n = 22), *E.*

255    *pseudokerneri* (n = 9), and *E. rostkoviana* (n = 9) had GS ranges greater than 1.1-fold. While

256    individuals with different GS values were often found in distant populations, such as in *E.*

257    *anglica* (0.6 pg and 0.73 pg, 525 km apart), and in *E. arctica* (1.04 pg and 1.33 pg, 903 km

258    apart), we also found considerable GS variation between populations in close proximity in *E.*

259    *foulaensis* (0.99 pg and 1.25 pg, 2.5 km apart) and in *E. confusa* (1.14 pg and 1.32 pg, same

260    population). In all cases, tests to distinguish genuine intraspecific variation from technical

261    artefacts confirmed the GS differences reported between individuals (see Methods and

262    Supporting Information Figure S**1c** and **d**). Generally, we found wider GS ranges in taxa with

263    more populations sampled. A notable exception was *E. micrantha* (GS range 1.14-1.21 pg from

264    17 individuals analysed from 9 populations, up to 962 km apart), which is discussed below.

265    In ANOVAs, the vast majority of the overall GS variation was explained by 'ploidy', while 'taxon'

266    and 'population' accounted for smaller significant fractions (Table **1**). 'Population' accounted for

267    considerably more variation than 'taxon' – 3 or 8 times, depending on whether hybrids were

268    included in the analysis or not. This difference is due to the few data available for most hybrids

269    (Figure **2a**, Supplementary Information Table **S1**). The fact that 'taxon' generally accounts for

270    only a small amount of variance is reflected by the near-continuous distribution of GSs within

271    each ploidy level (Figure **2b**). The distribution of tetraploid GS values has two gaps, caused by a

10

272    few exceptional individuals with extreme outliers in their GS values. While most tetraploid GS

273    values are between 1.07 and 1.26 pg (red horizontal lines in Figure **2b**), six samples had lower

274    (*E. arctica*, *E. foulaensis*, and *E. foulaensis* x *marshallii*), and seven higher, GS (*E. arctica*).

275

276    **Table 1. Partitioning of GS variation across *Euphrasia* species and hybrids.** Top analysis

277    includes all 192 samples from 13 species and 10 hybrids, and the lower analysis 157 samples

278    comprising just the 13 species (lower analysis). Both ANOVA tables detail the variance

279    components (Sum Sq) accounted for by ploidy, taxon and population.

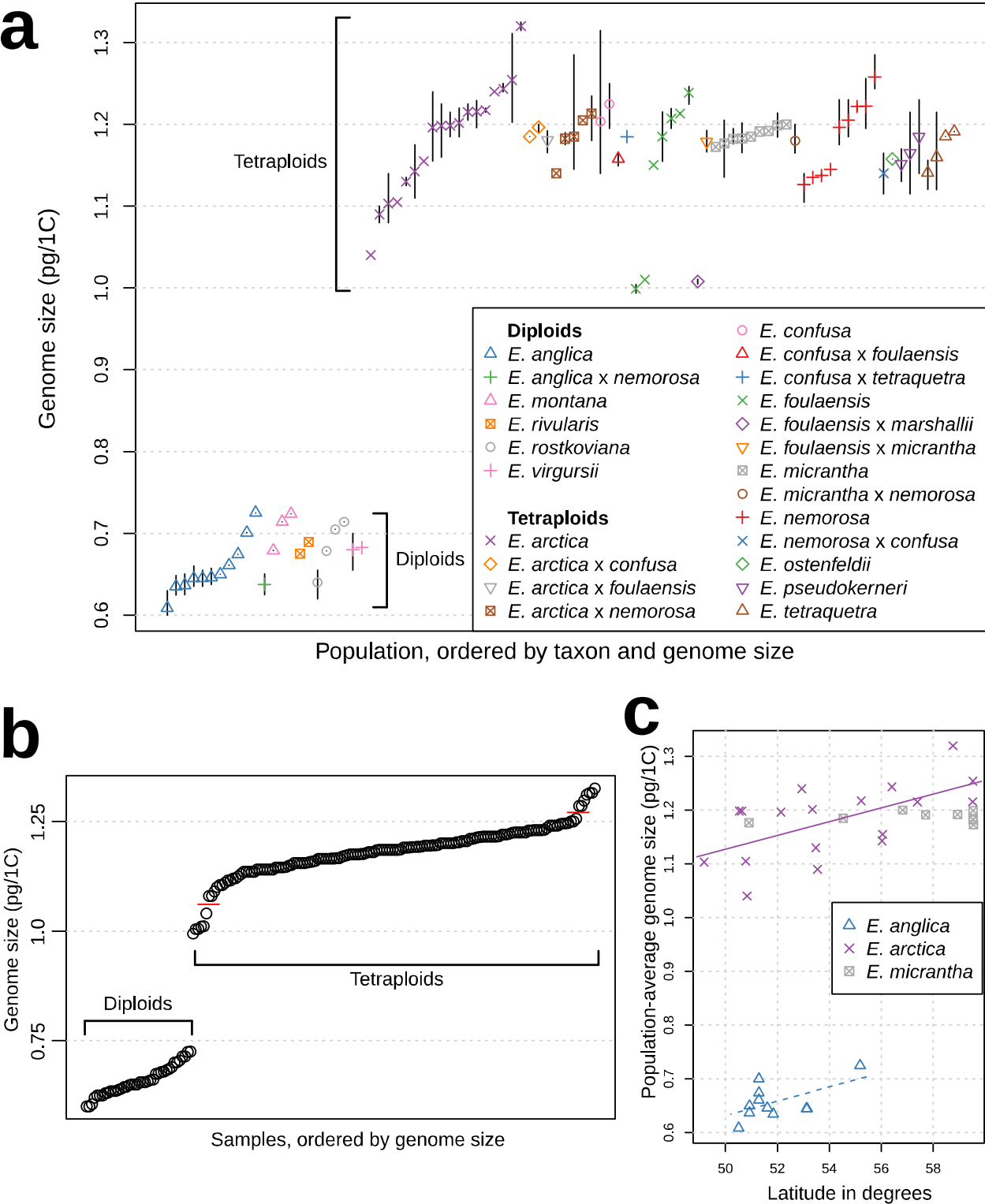|  |  | df | Sum Sq | Mean Sq | $F$ | $p$ |
|---|---|---|---|---|---|---|
| With hybrids | Ploidy | 1 | 8.67 | 8.67 | 9505.96 | **$< 2.0 \times 10^{-16}$** |
|  | Taxon | 21 | 0.11 | 0.01 | 6.00 | **$4.1 \times 10^{-10}$** |
|  | Population | 67 | 0.34 | 0.01 | 5.48 | **$7.8 \times 10^{-15}$** |
|  | Residuals | 102 | 0.09 | 0.00 |  |  |
| Without hybrids | Ploidy | 1 | 7.96 | 7.96 | 8763.74 | **$< 2.0 \times 10^{-16}$** |
|  | Taxon | 11 | 0.04 | 0.00 | 4.17 | **$6.9 \times 10^{-5}$** |
|  | Population | 62 | 0.33 | 0.01 | 5.92 | **$1.5 \times 10^{-13}$** |
|  | Residuals | 82 | 0.07 | 0.00 |  |  |

280

**Figure 2. Patterns of GS variation in British *Euphrasia*. a** The distribution of population-average GS for 90 populations of 23 taxa (13 species and 10 hybrids). Vertical bars indicate the GS range within each population where more than one individual was analysed. **b** Distribution of individual GS estimates for all 192 samples. Horizontal red lines indicate the limits of the

286    continuous part of the tetraploid GS distribution. **c** Population average genome sizes plotted

287    against latitude for the three most widely sampled species. The solid purple line indicates a

288    significant statistical relationship of GS with latitude across 17 populations of *E. arctica*. This

289    relationship was only marginally significant for 11 populations of *E. anglica* (dashed blue line).

290    No significant association was found across nine populations of the highly selfing *E. micrantha*.

291

292

293    Analyses of the three geographically widespread species with wider population sampling

294    revealed that GS variation was significantly partitioned by population for mixed-mating *E.*

295    *anglica* ($F_{10,12}$=9.86, $p$=2.3×10$^{-4}$) and *E. arctica* ($F_{17,25}$=10.5, $p$ < 1.7×10$^{-7}$), but not for highly

296    selfing *E. micrantha* ($F_{8,8}$=0.31, $p$=0.94). Further, the variance in population average GS was

297    significantly lower in *E. micrantha* than in *E. anglica* ($F_{10,8}$=11.65, $p$=9.6×10$^{-4}$) or *E. arctica*

298    ($F_{17,8}$=53. 2, $p$=2.3×10$^{-6}$).

299    Individual-based Mantel tests to link geographic distance and GS variation were significant over

300    all 40 diploid samples (Mantel statistic $r$=0.25, $p$=0.001) and all 152 tetraploids ($r$=0.04, $p$=0.01).

301    We then carried out Mantel tests based on population averages to exclude the very local

302    distance component. These tests were significant over all diploids ($r$=0.27, $p$=0.002) but not

303    over all tetraploid populations ($r$=0.04, $p$=0.09). However, *E. arctica*, the most widespread

304    tetraploid species, showed a pattern of isolation-by-distance at this level ($r$=0.24, $p$=0.015).
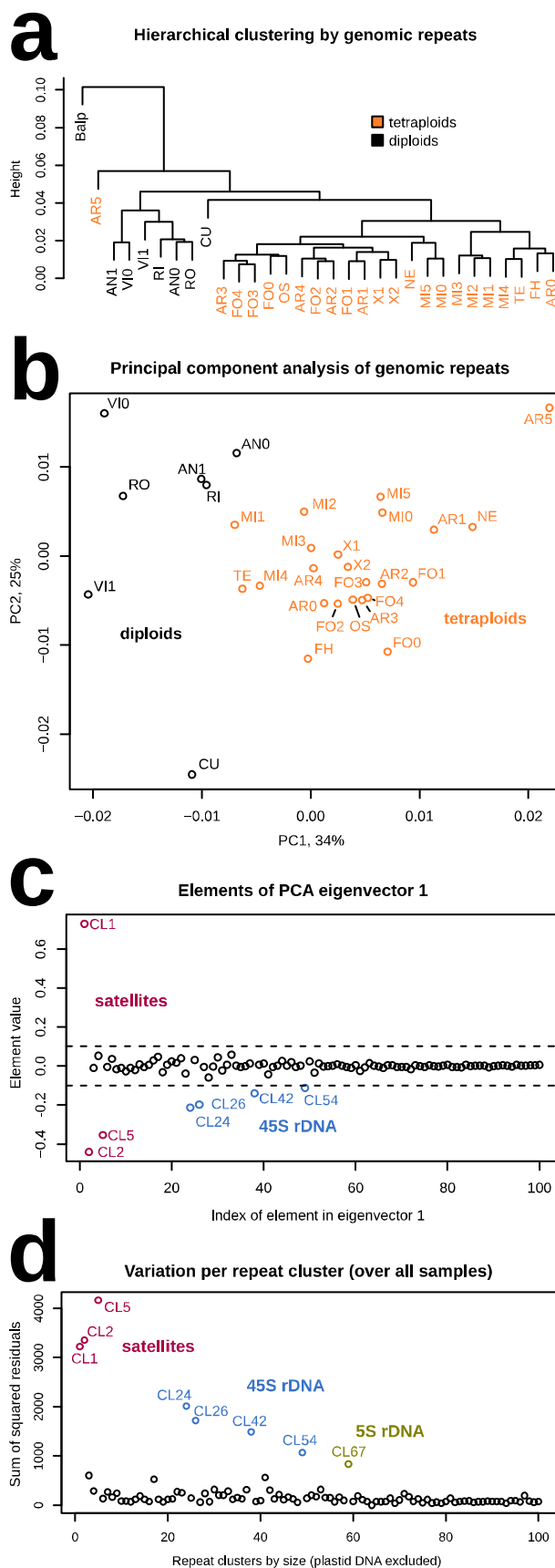
305    We confirmed a strong relationship between ploidy and latitude (ANOVA $F_{1,190}$=18.79,

306    $p$=2.4×10$^{-5}$), with diploids generally limited to lower latitudes (being particularly abundant in

307    southern England, Supporting Figure S2) while tetraploids extend to the very north of Britain.

308    However, there was no significant association between GS and latitude within ploidy levels

309    (treating taxon as a random effect, $t$=0.63, $p$=0.53). We then analysed the data for each of the

310    three widely sampled species individually using linear models (Figure **2c**). There was a non-

311    significant trend for the diploid *E. anglica* (slope=0.013pg/(degree latitude), $F_{1,9}$=4.23, $p$=0.07,

312    $r^2$=0.24). Of the tetraploids, *E. arctica* did show that GS increases significantly with latitude

313    (slope=0.013 pg/(degree latitude), $F_{1,16}$=9.36, $p$=0.008, $r^2$=0.31), whereas *E. micrantha* did not

314    ($F_{1,7}$=0.34, $p$=0.577).

315

316    *Variation in genomic repeat content*

13

317    To investigate the nature of the PAVs underpinning GS variation, we analysed the genomic

318    repeat content from whole genome sequencing data in 31 samples using the RE pipeline. RE's

319    output includes a set of annotated repeat clusters, representing individual repeat types. Our

320    samples included *B. alpina* (Orobanchaceae), 29 British *Euphrasia* samples (six diploids and 23

321    tetraploids), and one Austrian diploid (Supporting Information Table **S2**). Overall, 69.9% of all

322    *Euphrasia* reads analysed were identified as derived from repetitive DNA (i.e. they formed

323    repeat clusters with genome proportions > 0.01%). The average genomic repeat contents of

324    diploid and tetraploid *Euphrasia* samples differed, being 71.4% and 69.1%, respectively

325    ($F_{1,28}$=8.14, *p*=0.008). The repeat content for *B. alpina* was only 42.4%, which is an under-

326    estimate because repeats private to the species may have failed to form individual clusters

327    given our sampling design and cut-off threshold.

328    The most abundant repeat family, ranging from 25% in *E. anglica* (AN1) to 30% in *E. cuspidata*

329    (CU), was Angela, a type of Ty1/Copia long terminal repeat retrotransposon (LTR), which is

330    typically c. 8.5 kbp in length. Overall, all types of Ty1/Copia elements identified accounted for

331    30-39% of each *Euphrasia* genome, while Ty3/Gypsy elements typically occupied just 3-6% of

332    the genome (Supplementary Information Table S2).

333    To assess how well genomic repeat profiles in samples from different populations correspond

334    with species identity based on morphology, we used two unsupervised machine learning

335    techniques: hierarchical clustering and principal component analysis (PCA). We focussed our

336    analyses on the largest 100 repeat clusters, which together account for approximately 50% of

337    each genome, no matter if diploid or tetraploid. Each smaller repeat cluster had a genomic

338    proportion of < 0.7% in each sample. Hierarchical clustering resulted in a tree that grouped

339    samples largely by ploidy, rather than species identity, with the exception of (i) a sample of the

340    Austrian alpine *E. cuspidata* (CU), a species considered diploid, which grouped as sister to the

341    tetraploids, and (ii) tetraploid *E. arctica* from Cornwall (AR5), which grouped as sister to all other

342    *Euphrasia* samples (Figure **3a**). All species with multiple samples formed mixed branches with

343    other species in this tree. Among  the sympatric samples from Fair Isle, *E. micrantha* (MI1-3)

344    clustered separately from *E. arctica* (AR1-3) and *E. foulaensis* (FO1-4), both of which were

345    mixed with other species, similar to previous patterns of clustering from SNP-based analyses

346    (Becher *et al.*, 2020).

14

## a Hierarchical clustering by genomic repeats



## b Principal component analysis of genomic repeats



## c Elements of PCA eigenvector 1



## d Variation per repeat cluster (over all samples)



347

15

348 **[one col] Figure 3. Clustering of *Euphrasia* samples based on genomic repeat content. a**
349 Hierarchical clustering shows grouping largely by ploidy. **b** A PCA of the relative proportions of
350 the top 100 repeat clusters in 30 samples of *Euphrasia*. Diploids are shown in black and
351 tetraploids in orange. **c** Contribution of each repeat cluster to the first principal component (of
352 panel b). Clusters with negative values are enriched in diploids while those with positive values
353 are enriched in tetraploids. **d** The extent of variation in the genomic proportions of all individuals
354 for each repeat cluster. The codes in **a** and **b** are: Balp-*Bartsia alpina* (outgroup), five diploid
355 species (seven samples): AN-*E. anglica*, CU-*E. cuspidata*, VI-*E. vigursii*, RI-*E. rivularis*, RO-*E.*
356 *rostkoviana*, seven tetraploid species and two tetraploid hybrids: AR-*E. arctica*, FO-*E.*
357 *foulaensis*, MI-*E. micrantha*, NE-*E. nemorosa*, FH-*E. fharaidensis*, OS-*E. ostenfeldii*, TE-*E.*
358 *tetraquetra*, and X-tetraploid hybrids.

359 PCA without the outgroup *B. alpina* yielded a PC1 that explained 34% of the variance in our
360 repeat data, separating the diploid and tetraploid samples (Figure **3b**), whereas there was no
361 clear separation by species. The samples for some species were spread widely across the plot
362 (e.g. *E. arctica* (AR0-5) and *E. vigursii* (VI0, VI1)), while those of *E. micrantha* (MI0-5) grouped
363 relatively tightly. Although this does not preclude the possibility of species-specific repeat
364 patterns in *Euphrasia*, it is clear that there are no major differences in the relative abundance of
365 the common repeat types between the species. Within the 138 largest repeat clusters, none
366 was species-specific (i.e. present in individuals of only one species). Within the largest 701
367 clusters, none was diagnostic for a species (i.e., none was present in all samples of one species
368 but absent in all other samples).

369 To further analyse which repeat clusters separate diploids and tetraploids in the PCA (Figure
370 **3b**), we plotted the elements of eigenvector 1, which correspond to the effect of each repeat
371 cluster on the position of a sample along PC1 (Figure **3c**). Seven repeat clusters have a large
372 effect on PC1, the satellite clusters CL1, CL2 and CL5, and all clusters of the 45S ribosomal
373 DNA (CL24, CL26, CL42, and CL56). Satellite clusters CL1 and CL2 have monomer size peaks
374 of approximately 145 nucleotides as commonly seen in centromeric repeats. In addition, some
375 reads of CL1 and CL2 had paired-end mates in CL22, indicating physical proximity of the
376 repeats within the genome. CL22, in turn, had been annotated as CRM, which is a type of
377 Ty3/Gypsy chromovirus retrotransposon that commonly targets centromeric sequences (Nagaki
378 *et al.*, 2003; Neumann *et al.*, 2011).

379 Among all 17 broad repeat types identified by RE (see Supplementary Information Table S2),
380 we found significant differences between ploidy levels for two. Diploid genomes contained
381 higher average proportions of 45S rDNA (4.9%) than tetraploids (2.0%, $F_{1, 28}=20.4$, $p_{corr}<0.001$),
382 with the genomic proportion ranging from 1.7% to 5.7% in diploids and from 0.8% to 3.4% in

16

383    tetraploids. Tetraploids contained, on average, more Ty1/copia Ale elements (0.15%) than

384    diploids (0.09%, $F_{1,28}$=11.18, $p_{corr}$=0.018). While our PCA approach had identified some

385    satellites as highly differentiated in copy number (see above), differences over all satellites were

386    not significant. This is because there was differential enrichment in the ploidy levels for CL1

387    versus CL2 and CL5 (Figure **3c**). Overall, there is comparatively little differentiation in genomic

388    repeats between the ploidy levels.

389    We also assessed the variation in repeat content over all samples for each repeat cluster. The

390    eight most variable clusters (i.e. having the biggest differences in repeat proportions between

391    individuals, Figure **3d**), are all tandem repeats (satellites including rDNA). The first seven are

392    the same repeats that separated the ploidy levels in the PCA. The eighth most variable repeat

393    (CL67), which is variable in both ploidy levels, corresponds to the 5S rDNA.

394    Of the samples analysed with RE, nine tetraploids were from populations which also had GS

395    estimates obtained in this study. Testing the largest 100, 200, and 1000 repeat clusters for

396    correlations between GS and abundance of individual repeat clusters, and correcting for

397    multiple testing by Bonferroni correction, no repeat cluster showed a significant correlation

398    between its abundance in an individual and the population-average genome size. All evidence

399    from repetitive elements suggests that the GS differences between *Euphrasia* individuals of the

400    same ploidy levels are not due to large changes in the genomic proportion of any one specific

401    repeat.

402

403    **Discussion**

404    In this study, we investigated the nature of GS variation across taxonomically complex diploid

405    and tetraploid British *Euphrasia*. We complemented an extensive population survey of GS

406    variation with an analysis of genomic repeat composition from seven diploids and 23 tetraploid

407    *Euphrasia*. Overall, we find notable GS variation between populations of the same species,

408    representing a wide range of genuine intraspecific GS variation. Within ploidy levels there is a

409    continuum of GS variation, though ploidy levels have discrete GS ranges. These differences

410    within and between ploidy levels are not attributable to large copy number changes of an

411    individual DNA repeat, but rather to multiple segregating PAVs. Here, we first discuss the link

412    between GS variation and population dynamics and speciation history, highlighting how GS is

413    shaped by many similar processes as population-level sequence variation. We then consider

17

414    the landscape of repeat dynamics and the potential association with *Euphrasia* polyploid

415    genome history. Finally, we consider the wider implications of framing GS variation in a

416    population genetic framework.

417    *Genome size variation mirrors population genetic patterns*

418    *Euphrasia* are renowned as a taxonomically complex group where species are recent in origin

419    and show subtle morphological differences, and taxa readily hybridise in areas of secondary

420    contact (Gussarova *et al.*, 2008; Wang *et al.*, 2018a). Previous population genetic analysis have

421    shown genetic variation is not clearly partitioned by species (Kolseth & Lönn, 2005; French *et*

422    *al.*, 2008; Becher *et al.*, 2020), particularly in widespread co-occurring outcrossers, with only

423    certain taxa, like the moorland selfing species *E. micrantha*, being genetically diverged. Here,

424    we find GS variation mirrors these findings of population genetic structure inferred from

425    molecular data. Our results add doubt to the distinctiveness of species, with taxa clearly not

426    showing distinct GS ranges indicative of reproductive isolation. Moreover, previous findings

427    have reported a considerably higher mean GS of 2.73 pg for five samples of diploid *E.*

428    *rostkoviana* from Bosnia and Herzegovina (Siljak-Yakovlev *et al.*, 2010) compared with our

429    estimates that ranged 0.62 -0.71 pg. This notable discrepancy raises a number of non-mutually

430    exclusive hypotheses: (1) heterogeneous GS variation within currently named species may be a

431    consequence of different taxonomic concepts applied across Europe; (2) lower GS variation

432    within British *Euphrasia* may be a consequence of hybridisation and homogenisation of GS

433    variation in Britain or a distinct polyploid history elsewhere in Europe; (3) identification problems

434    or technical issues may affect previous GS estimates.

435    The continuous GS distribution across species boundaries within ploidy levels in *Euphrasia*

436    resembles the findings of Hanušová *et al.* (2014) for species of the lycophyte *Diphasiastrum* at

437    allopatric and sympatric sites. These authors concluded that considerable GS variation within

438    species resulted from introgression from other sympatric species. Depending on the sizes and

439    number of segregating PAVs (see Figure **1b** and **c**), hybridisation between divergent

440    populations may homogenise local GS, or introduce GS differences. In our study, three

441    populations from Fair Isle (one *E. foulaensis* x *E. marshallii* and two *E. foulaensis*) located within

442    5 km of each other show likely signals of introgression. Their GS estimates were more than 10%

443    lower than the mean GS of all tetraploids, including all other Fair Isle samples (Figure **2a**). While

444    these populations might have independently evolved lower GS, it seems more plausible that

445    they share large GS difference variants (such as missing dispensable chromosomes or

18

446    chromosome regions, Figure **1c**). An explanation of genomic homogenisation in sympatry is in

447    keeping with the growing body of plant research showing gene flow at the early stages of

448    species divergence, or between closely related species (e.g. Strasburg & Rieseberg, 2008;

449    Papadopulos *et al.*, 2011; Brandvain *et al.*, 2014; Sawangproh *et al.*, 2020). Such observations

450    of divergence with gene flow are often coupled with species differences being maintained by a

451    few diverged regions under strong selection maintaining species identities (e.g. Twyford &

452    Friedman, 2015), a possibility we are currently investigating in *Euphrasia*.

453    Within three of the widespread species that we sampled extensively, we found considerably

454    higher GS variation in the mainly outcrossing *E. anglica* and *E. arctica* than in highly selfing *E.*

455    *micrantha*. Unlike the outcrossing species, *E. micrantha* shows no increase in GS at higher

456    latitudes. Lower diversity is expected for several reasons in young selfing lineages such as *E.*

457    *micrantha*. Firstly, selfing reduces the effective population size, resulting in lower genetic

458    variation (Nordborg, 1997), presumably including PAVs. Secondly, the reduced effective rate of

459    crossing over between the chromosomes of a selfing species further reduces the effective

460    population size (Conway *et al.*, 1999). Thirdly, selfing species are rarely polymorphic for B

461    chromosomes (Burt & Trivers, 2008), one source of GS variation in the Orobanchaceae, for

462    instance in closely related *Rhinanthus* (Wulff, 1939; Hambler, 1953). Finally, partially selfing

463    species are less likely to acquire GS variants through introgression (e.g. Pajkovic *et al.*, 2014).

464    Older highly selfing lineages may, however, have diversified ecologically and become restricted

465    to different habitats, and might evolve GS differences.

466

467    *Genome size differences and genomic repeats*

468    We found very low differentiation of genomic repeats between species of British *Euphrasia*, with

469    few species-specific repeats. Consistent with phylogenetic work (Gussarova *et al.*, 2008; Wang

470    *et al.*, 2018a), there were no examples where all species samples cluster together based on

471    repeat content (Figure **3a**). The fact that species of British *Euphrasia* are closely related and

472    often hybridise, makes lineage-specific large-scale gains or losses of individual repeat groups,

473    as seen in other plants (Piegu *et al.*, 2006; Macas *et al.*, 2015; McCann *et al.*, 2020), an unlikely

474    cause for the observed GS variation. Instead, the observed differences are likely due to

475    changes in numerous different repeats segregating within the *Euphrasia* gene pool. At present,

476    it is hard to tell whether these PAVs comprise numerous individual repeat copies or whether

19

477    there are (also) larger-scale PAVs like the loss or gain of chromosome fragments as

478    hypothesised in hybridising species of *Anacyclus* (Agudo *et al.*, 2019; Vitales *et al.*, 2020). The

479    high frequency of hybridisation in *Euphrasia* may lead to increased levels of structural

480    rearrangements due to ectopic recombination, which may be more common between

481    heterozygous genomic repeats (Morgan, 2001).

482    Between ploidy levels of *Euphrasia*, we found allotetraploids had an 11% lower mean GS

483    compared with the value predicted from doubling the mean GS of diploids. This discrepancy

484    may have originated from genome downsizing, commonly seen during re-diploidisation. It may

485    also be explained by the fusion of two diverged diploid genomes of different size, as seen in

486    allopolyploid *Gossypium* (Hendrix & Stewart, 2005) and *Arabidopsis suecica* (Burns *et al.*,

487    2021). However, the absence of interploidy repeat divergence in *Euphrasia* differs from other

488    allotetraploid systems, where diverged sub-genomes tend to show differences in genomic

489    repeats (Zhao *et al.*, 1998; Hawkins *et al.*, 2006; Renny-Byfield *et al.*, 2015; Dodsworth *et al.*,

490    2020). This lack of repeat differentiation is notable because nuclear k-mer spectra (Becher *et*

491    *al.*, 2020) and rDNA sequences (Wang *et al.*, 2018a) suggest considerable sequence

492    divergence between the tetraploid sub-genomes, corresponding to a split of approximately 8

493    million years (Gussarova *et al.*, 2008).

494    Tandem repeats such as rDNA and other satellite DNAs are generally found to be the fastest

495    evolving fraction of the repeatome, showing divergence in both copy number and sequence

496    between closely related species (e.g. Tek *et al.*, 2005; Ambrozová *et al.*, 2011; Renny-Byfield *et*

497    *al.*, 2012; Becher *et al.*, 2014; Ávila Robledillo *et al.*, 2020) and populations (Ananiev *et al.*,

498    1998). We confirmed this in *Euphrasia*, where tandem repeats accounted for the eight repeat

499    clusters with the highest inter-individual variation in genomic abundance (Figure **3d**). While

500    differing across individuals, repeat content did not show any clear signal of divergence between

501    particular species. For example, the comparison between *E. micrantha* and divergent tetraploids

502    such as *E. arctica*, did not reveal a signal of divergence in repeat content. This is surprising not

503    just because of their morphological distinctiveness, but their difference in outcrossing rate, with

504    theory predicting that the copy-number and equilibrium frequency of transposable elements

505    depends on the level of selfing in a population (Morgan, 2001; Dolgin & Charlesworth, 2006). A

506    likely explanation is that the shift to high-selfing in *E. micrantha* is relatively recent compared to

507    the time it takes for the genomic repeat content to reach equilibrium level.

508

509   *Evolution of genome size variation*

510   The continuous GS variation within and between *Euphrasia* species, coupled with these

511   differences likely being a product of segregating PAV across the genome, underlines the

512   polygenic nature of GS variation. Regarding GS differences as the result of segregating (i.e.

513   genetic) variants blurs the classic distinction between genotype and nucleotype, where

514   "nucleotype" refers to "conditions of the nucleus that affect the phenotype independently of the

515   informational content of the DNA", essentially identical to GS (Bennett, 1971, 1977). Because

516   GS has been shown to be correlated with many traits including cell size, stomatal pore size, the

517   duration of cell division, and life-history differences (e.g. Šímová & Herben, 2012; Bilinski *et al.*,

518   2018; Roddy *et al.*, 2020), it is plausible the GS is affected indirectly by selection on such traits.

519   There might be additional indirect selection on GS according according to the mutational-hazard

520   hypothesis (e.g. Lynch, 2011), which proposes that large GS may be selected against because

521   there is more opportunity for the accumulation of deleterious mutations.

522

523   It follows that individual PAVs may be under different kinds of simultaneous selection, potentially

524   of different directionality. For instance, there might be positive selection on an adaptive

525   insertion, which is simultaneously selected against because it increases GS. Further, because

526   selection at one locus affects regions that are physically linked (i.e. selection at linked sites,

527   Maynard Smith & Haigh, 1974; Charlesworth et al., 1993), the footprint of selection on genome

528   regions is modified by the (effective) rate of crossing over, which varies along genomes and

529   between mating systems.

530

531   Research on GS is somewhat decoupled from studies on sequence-based variation in

532   populations. We suggest future research into GS evolution should consider both patterns of total

533   GS and the population processes underlying this variation. In addition to furthering our

534   understanding of intraspecific GS diversity in *Euphrasia* and other plant groups, answers to

535   these questions will also improve our understanding of GS evolution between species and

536   across phylogenies, which starts at the population level.

537

21

550
551
552

## Author Contributions

554 • HB analysed the data with input from MRB and ADT

555 • ADT, CM, HB, and MRB collected samples.

556 • CM confirmed species identifications.

557 • ADT and IJL designed the study.

558 • RFP, JP, and IJL generated the GS data.

559 • HB and ADT wrote the manuscript.

560 • All authors read and commented on the manuscript.

561

## Data Availability

563 The whole genome-sequencing data are available from the sequence read archive, Bioprojects
564 PRJNA624746 and PRJNA678958. The scripts and data required to replicate our results are
565 available from GitHub, repository: zzzzzzz (to be added upon acceptance).
566

## References

567

568 **Achigan-Dako EG, Fuchs J, Ahanchede A, Blattner FR**. **2008**. Flow cytometric analysis in

569 *Lagenaria siceraria* (Cucurbitaceae) indicates correlation of genome size with usage types and

570 growing elevation. *Plant Systematics and Evolution* **276**: 9.

571 **Agudo AB, Torices R, Loureiro J, Castro S, Castro M, Álvarez I**. **2019**. Genome Size

572 Variation in a Hybridizing Diploid Species Complex in *Anacyclus* (Asteraceae: Anthemideae).

573 *International Journal of Plant Sciences* **180**: 374–385.

574 **Ambrozová K, Mandáková T, Bures P, Neumann P, Leitch IJ, Koblízková A, Macas J,**

575 **Lysak M a**. **2011**. Diverse retrotransposon families and an AT-rich satellite DNA revealed in

576 giant genomes of *Fritillaria* lilies. *Annals of Botany* **107**: 255–268.

577 **Ananiev E V, Phillips RL, Rines HW**. **1998**. A knob-associated tandem repeat in maize

578 capable of forming fold-back DNA segments: Are chromosome knobs megatransposons?

579 *Proceedings of the National Academy of Sciences* **95**: 10785 LP – 10790.

580 **Ávila Robledillo L, Neumann P, Koblížková A, Novák P, Vrbová I, Macas J**. **2020**.

581 Extraordinary sequence diversity and promiscuity of centromeric satellites in the legume tribe

582 Fabeae. *Molecular Biology and Evolution* **37**: 2341–2356.

583 **Bainard JD, Newmaster SG, Budke JM**. **2019**. Genome size and endopolyploidy evolution

584 across the moss phylogeny. *Annals of Botany* **125**: 543–555.

585 **Beaulieu JM, Leitch IJ, Patel S, Pendharkar A, Knight CA**. **2008**. Genome size is a strong

586 predictor of cell size and stomatal density in angiosperms. *New Phytologist* **179**: 975–986.

587 **Becher H, Brown MR, Powell G, Metherell C, Riddiford NJ, Twyford AD**. **2020**. Maintenance

588 of species differences in closely related tetraploid parasitic *Euphrasia* (Orobanchaceae) on an

589 isolated island. *Plant Communications*: 100105.

590 **Becher H, Ma L, Kelly LJ, Kovařík A, Leitch IJ, Leitch AR**. **2014**. Endogenous pararetrovirus

591 sequences associated with 24 nt small RNAs at the centromeres of *Fritillaria imperialis* L.

592 (Liliaceae), a species with a giant genome. *The Plant journal : for cell and molecular biology*

593 **80**: 823–833.

594 **Bennett MD**. **1971**. The duration of meiosis. *Proceedings of the Royal Society of London.*

595    *Series B. Biological Sciences* **178**: 277–299.

596    **Bennett MD**. **1977**. The time and duration of meiosis. *Philosophical Transactions of the Royal*
597    *Society of London. B, Biological Sciences* **277**: 201–226.

598    **Bennett MD, Smith JB**. **1991**. Nuclear DNA amounts in angiosperms. *Philosophical*
599    *Transactions of the Royal Society of London. B, Biological Sciences* **334**: 309–345.

600    **Bilinski P, Albert PS, Berg JJ, Birchler JA, Grote MN, Lorant A, Quezada J, Swarts K,**
601    **Yang J, Ross-Ibarra J**. **2018**. Parallel altitudinal clines reveal trends in adaptive evolution of
602    genome size in *Zea mays. PLOS Genetics* **14**: e1007162.

603    **Blommaert J**. **2020**. Genome size evolution: towards new model systems for old questions.
604    *Proceedings of the Royal Society B: Biological Sciences* **287**: 20201441.

605    **Brandvain Y, Kenney AM, Flagel L, Coop G, Sweigart AL**. **2014**. Speciation and
606    Introgression between Mimulus nasutus and Mimulus guttatus. *PLoS Genetics* **10**: e1004410.

607    **Brown MR, Frachon N, Wong ELY, Metherell C, Twyford AD**. **2020**. Life history evolution,
608    species differences, and phenotypic plasticity in hemiparasitic eyebrights (*Euphrasia*). *American*
609    *Journal of Botany* **107**: 456–465.

610    **Burns R, Mandáková T, Gunis J, Soto-Jiménez LM, Liu C, Lysak MA, Novikova PY,**
611    **Nordborg M**. **2021**. Gradual evolution of allopolyploidy in *Arabidopsis suecica*. *bioRxiv*:
612    2020.08.24.264432.

613    **Burt A, Trivers R**. **2008**. *Genes in conflict: The biology of selfish genetic elements*. Cambridge
614    (Massachusetts): Harvard University Press.

615    **Cacho NI, McIntyre PJ, Kliebenstein DJ, Strauss SY**. **2021**. Genome size evolution is
616    associated with climate seasonality and glucosinolates, but not life history, soil nutrients or
617    range size, across a clade of mustards. *Annals of Botany.*

618    **Chia J-M, Song C, Bradbury PJ, Costich D, de Leon N, Doebley J, Elshire RJ, Gaut B,**
619    **Geller L, Glaubitz JC, *et al.* 2012**. Maize HapMap2 identifies extant variation from a genome in
620    flux. *Nature Genetics* **44**: 803–807.

621    **Conway DJ, Roper C, Oduola AMJ, Arnot DE, Kremsner PG, Grobusch MP, Curtis CF,**

622 **Greenwood BM**. **1999**. High recombination rate in natural populations of *Plasmodium*

623 *falciparum*. *Proceedings of the National Academy of Sciences* **96**: 4506 LP – 4511.

624 **Costich DE, Meagher TR, Yurkow EJ**. **1991**. A rapid means of sex identification in *Silene*

625 *latifolia* by use of flow cytometry. *Plant Molecular Biology Reporter* **9**: 359–370.

626 **Díez CM, Gaut BS, Meca E, Scheinvar E, Montes-Hernandez S, Eguiarte LE, Tenaillon MI**.

627 **2013**. Genome size variation in wild and cultivated maize along altitudinal gradients. *New*

628 *Phytologist* **199**: 264–276.

629 **Dodsworth S, Guignard MS, Pérez-Escobar OA, Struebig M, Chase MW, Leitch AR**. **2020**.

630 Repetitive DNA restructuring across multiple *Nicotiana* allopolyploidisation events shows a lack

631 of strong cytoplasmic bias in influencing repeat turnover. *Genes* **11**.

632 **Doležel J, Bartoš J, Voglmayr H, Greilhuber J**. **2003**. Letter to the editor. *Cytometry* **51A**:

633 127–128.

634 **Dolgin ES, Charlesworth B**. **2006**. The fate of transposable elements in asexual populations.

635 *Genetics* **174**: 817–827.

636 **Duchoslav M, Šafářová L, Jandová M**. **2013**. Role of adaptive and non-adaptive mechanisms

637 forming complex patterns of genome size variation in six cytotypes of polyploid *Allium*

638 *oleraceum* (Amaryllidaceae) on a continental scale. *Annals of Botany* **111**: 419–431.

639 **Fisher RA**. **1919**. XV.—The correlation between relatives on the supposition of mendelian

640 inheritance. *Transactions of the Royal Society of Edinburgh* **52**: 399–433.

641 **French GC, Hollingsworth PM, Silverside AJ, Ennos RA**. **2008**. Genetics, taxonomy and the

642 conservation of British *Euphrasia*. *Conservation Genetics* **9**: 1547–1562.

643 **Gao L, Gonda I, Sun H, Ma Q, Bao K, Tieman DM, Burzynski-Chang EA, Fish TL,**

644 **Stromberg KA, Sacks GL, *et al.* 2019**. The tomato pan-genome uncovers new genes and a

645 rare allele regulating fruit flavor. *Nature Genetics* **51**: 1044–1051.

646 **Göktay M, Fulgione A, Hancock AM**. **2020**. A new catalog of structural variants in 1,301 *A.*

647 *thaliana* lines from africa, eurasia, and north america reveals a signature of balancing selection

648 at defense response genes. *Molecular Biology and Evolution*.

649  **Gregory TR, Johnston JS**. **2008**. Genome size diversity in the family Drosophilidae. *Heredity*
650  **101**: 228–238.

651  **Greilhuber J**. **2005**. Intraspecific variation in genome size in angiosperms: Identifying its
652  existence. *Annals of Botany* **95**: 91–98.

653  **Greilhuber J, Doležel J, Lysák MA, Bennett M**. **2005**. The origin, evolution and proposed
654  stabilization of the terms 'Genome Size' and 'C-Value' to describe nuclear DNA contents.
655  *Annals of Botany* **95**: 255–260.

656  **Greilhuber J, Leitch IJ**. **2013**. Genome size and the phenotype: Physical Structure, Behaviour
657  and Evolution of Plant Genomes. In: Greilhuber J, Dolezel J, Wendel JF, eds. Plant Genome
658  Diversity Volume 2. Vienna: Springer, 323–344.

659  **Greilhuber J, Temsch EM, Loureiro JCM**. **2007**. Nuclear DNA content measurement. *Flow*
660  *Cytometry with Plant Cells*: 67–101.

661  **Guignard MS, Crawley MJ, Kovalenko D, Nichols RA, Trimmer M, Leitch AR, Leitch IJ**.
662  **2019**. Interactions between plant genome size, nutrients and herbivory by rabbits, molluscs and
663  insects on a temperate grassland. *Proceedings of the Royal Society B: Biological Sciences* **286**:
664  20182619.

665  **Guignard MS, Nichols RA, Knell RJ, Macdonald A, Romila C-A, Trimmer M, Leitch IJ,**
666  **Leitch AR**. **2016**. Genome size and ploidy influence angiosperm species' biomass under
667  nitrogen and phosphorus limitation. *New Phytologist* **210**: 1195–1206.

668  **Gussarova G, Popp M, Vitek E, Brochmann C**. **2008**. Molecular phylogeny and biogeography
669  of the bipolar *Euphrasia* (Orobanchaceae): Recent radiations in an old genus. *Molecular*
670  *Phylogenetics and Evolution* **48**: 444–460.

671  **Haberer G, Kamal N, Bauer E, Gundlach H, Fischer I, Seidel MA, Spannagl M, Marcon C,**
672  **Ruban A, Urbany C, *et al.* 2020**. European maize genomes highlight intraspecies variation in
673  repeat and gene content. *Nature Genetics* **52**: 950–957.

674  **Hambler DJ**. **1953**. Prochromosomes and supernumerary chromosomes in *Rhinanthus minor*
675  Ehrh. *Nature* **172**: 629–630.

676  **Hanušová K, Ekrt L, Vít P, Kolář F, Urfus T**. **2014**. Continuous morphological variation

26

677    correlated with genome size indicates frequent introgressive hybridization among *Diphasiastrum*

678    species (Lycopodiaceae) in Central Europe. *PLOS ONE* **9**: e99552.

679    **Hawkins JS, Kim H, Nason JD, Wing RA, Wendel JF**. **2006**. Differential lineage-specific

680    amplification of transposable elements is responsible for genome size variation in *Gossypium*.

681    *Genome Research* **16**: 1252–1261.

682    **Hendrix B, Stewart JM**. **2005**. Estimation of the nuclear DNA content of *Gossypium* species.

683    *Annals of Botany* **95**: 789–797.

684    **Hirsch CN, Foerster JM, Johnson JM, Sekhon RS, Muttoni G, Vaillancourt B,**

685    **Peñagaricano F, Lindquist E, Pedraza MA, Barry K, *et al.* 2014**. Insights into the maize pan-

686    genome and pan-transcriptome. *The Plant Cell* **26**: 121–135.

687    **Hübner S, Bercovich N, Todesco M, Mandel JR, Odenheimer J, Ziegler E, Lee JS, Baute**

688    **GJ, Owens GL, Grassa CJ, *et al.* 2019**. Sunflower pan-genome analysis shows that

689    hybridization altered gene content and disease resistance. *Nature Plants* **5**: 54–62.

690    **Kolář F, Čertner M, Suda J, Schönswetter P, Husband BC**. **2017**. Mixed-ploidy species:

691    progress and opportunities in polyploid research. *Trends in Plant Science* **22**: 1041–1055.

692    **Kolseth A-K, Lönn M**. **2005**. Genetic structure of *Euphrasia stricta* on the Baltic island of

693    Gotland, Sweden. *Ecography* **28**: 443–452.

694    **Koornneef M, Hanhart CJ, van der Veen JH**. **1991**. A genetic and physiological analysis of

695    late flowering mutants in *Arabidopsis thaliana*. *Molecular and General Genetics MGG* **229**: 57–

696    66.

697    **Leitch IJ, Beaulieu JM, Cheung K, Hanson L, Lysak M a, Fay MF**. **2007**. Punctuated genome

698    size evolution in Liliaceae. *Journal of evolutionary biology* **20**: 2296–308.

699    **Leitch IJ, Hanson L, Lim KY, Kovarik A, Chase MW, Clarkson JJ, Leitch AR**. **2008**. The ups

700    and downs of genome size evolution in polyploid species of Nicotiana (Solanaceae). *Annals of*

701    *Botany* **101**: 805–814.

702    **Leitch AR, Leitch IJ**. **2008**. Genomic plasticity and the diversity of polyploid plants. *Science*

703    **320**: 481–3.

704 **Long Q, Rabanal FA, Meng D, Huber CD, Farlow A, Platzer A, Zhang Q, Vilhjálmsson BJ,**

705 **Korte A, Nizhynska V,** *et al.* **2013**. Massive genomic variation and strong selection in

706 *Arabidopsis thaliana* lines from Sweden. *Nature Genetics* **45**: 884–890.

707 **Loureiro J, Rodriguez E, Doležel J, Santos C**. **2007**. Two new nuclear isolation buffers for

708 plant DNA flow cytometry: A test with 37 Species. *Annals of Botany* **100**: 875–888.

709 **Lynch M**. **2011**. Statistical inference on the mechanisms of genome evolution. *PLOS Genetics*

710 **7**: e1001389.

711 **Macas J, Novák P, Pellicer J, Čížková J, Koblížková A, Neumann P, Fuková I, Doležel J,**

712 **Kelly LJ, Leitch IJ**. **2015**. In depth characterization of repetitive DNA in 23 plant genomes

713 reveals sources of genome size variation in the legume tribe Fabeae. *PLOS ONE* **10**:

714 e0143424.

715 **McCann J, Macas J, Novák P, Stuessy TF, Villaseñor JL, Weiss-Schneeweiss H**. **2020**.

716 Differential genome size and repetitive DNA evolution in diploid species of Melampodium sect.

717 Melampodium (Asteraceae). *Frontiers in Plant Science* **11**: 362.

718 **Meagher TR, Gilies ACM, Costich DE**. **2005**. Genome size, quantitative genetics and the

719 genomic basis for flower size evolution in *Silene latifolia*. *Annals of Botany* **95**: 247–254.

720 **Metherell C, Rumsey FJ**. **2018**. *Eyebrights (*Euphrasia*) of the UK and Ireland* (J Edmondson,

721 Ed.). Bristol: Botanical Society of Britain and Ireland.

722 **Morgan MT**. **2001**. Transposable element number in mixed mating populations. *Genetical*

723 *Research* **77**: 261–275.

724 **Nagaki K, Song J**, **Stupar RM, Parokonny AS, Yuan Q, Ouyang S, Liu J, Hsiao J, Jones**

725 **KM, Dawe RK,** *et al.* **2003**. Molecular and cytological analyses of large tracks of centromeric

726 DNA reveal the structure and evolutionary dynamics of maize centromeres. *Genetics* **163**: 759

727 LP – 770.

728 **Napp-Zinn K**. **1955**. Genetische Grundlagen des Kältebedürfnisses bei *Arabidopsis thaliana*

729 (L.)Heynh. *Naturwissenschaften* **42**: 650.

730 **Neumann P, Navrátilová A, Koblížková A, Kejnovský E, Hřibová E, Hobza R, Widmer A,**

731 **Doležel J, Macas J**. **2011**. Plant centromeric retrotransposons: a structural and cytogenetic

732    perspective. *Mobile DNA* **2**: 4.

733    **Nordborg M**. **1997**. Structured coalescent processes on different time scales. *Genetics* **146**:

734    1501 LP – 1514.

735    **Novák P, Guignard MS, Neumann P, Kelly LJ, Mlinarec J, Koblížková A, Dodsworth S,**

736    **Kovařík A, Pellicer J, Wang W,** *et al.* **2020a**. Repeat-sequence turnover shifts fundamentally

737    in species with large genomes. *Nature Plants* **6**: 1325–1329.

738    **Novák P, Neumann P, Macas J**. **2010**. Graph-based clustering and characterization of

739    repetitive sequences in next-generation sequencing data. *BMC Bioinformatics* **11**: 378.

740    **Novák P, Neumann P, Macas J**. **2020b**. Global analysis of repetitive DNA from unassembled

741    sequence reads using RepeatExplorer2. *Nature Protocols* **15**: 3745–3776.

742    **Novák P, Neumann P, Pech J, Steinhaisl J, Macas J**. **2013**. RepeatExplorer: a Galaxy-based

743    web server for genome-wide characterization of eukaryotic repetitive elements from next-

744    generation sequence reads. *Bioinformatics* **29**: 792–793.

745    **Pajkovic M, Lappe S, Barman R, Parisod C, Neuenschwander S, Goudet J, Alvarez N,**

746    **Guadagnuolo R, Felber F, Arrigo N**. **2014**. Wheat alleles introgress into selfing wild relatives:

747    empirical estimates from approximate Bayesian computation in *Aegilops triuncialis*. *Molecular*

748    *Ecology* **23**: 5089–5101.

749    **Papadopulos AST, Baker WJ, Crayn D, Butlin RK, Kynast RG, Hutton I, Savolainen V**.

750    **2011**. Speciation with gene flow on Lord Howe Island. *Proceedings of the National Academy of*

751    *Sciences* **108**: 13188–13193.

752    **Pellicer J, Garcia S, Canela MÁ, Garnatje T, Korobkov AA, Twibell JD, Vallès J**. **2010**.

753    Genome size dynamics in *Artemisia* L. (Asteraceae): following the track of polyploidy. *Plant*

754    *Biology* **12**: 820–830.

755    **Pellicer J, Hidalgo O, Dodsworth S, Leitch IJ**. **2018**. Genome size diversity and its impact on

756    the evolution of land plants. *Genes* **9**: 88.

757    **Pellicer J, Leitch IJ**. **2020**. The Plant DNA C-values database (release 7.1): an updated online

758    repository of plant genome size data for comparative studies. *New Phytologist* **226**: 301–305.

759 **Pellicer J, Powell RF, Leitch IJ**. **2021**. The application of flow cytometry for estimating genome

760 size, ploidy level endopolyploidy, and reproductive modes in plants BT - Molecular Plant

761 Taxonomy: Methods and Protocols. In: Besse P, ed. New York, NY: Springer US, 325–361.

762 **Piegu B, Guyot R, Picault N, Roulin A, Sanyal A, Saniyal A, Kim H, Collura K, Brar DS,**

763 **Jackson S,** *et al.* **2006**. Doubling genome size without polyploidization: dynamics of

764 retrotransposition-driven genomic expansions in Oryza australiensis, a wild relative of rice.

765 *Genome research* **16**: 1262–9.

766 **R Core Team**. **2019**. R: A Language and Environment for Statistical Computing.

767 **Renner SS, Heinrichs J, Sousa A**. **2017**. The sex chromosomes of bryophytes: Recent

768 insights, open questions, and reinvestigations of *Frullania dilatata* and *Plagiochila asplenioides*.

769 *Journal of Systematics and Evolution* **55**: 333–339.

770 **Renny-Byfield S, Gong L, Gallagher JP, Wendel JF**. **2015**. Persistence of subgenomes in

771 paleopolyploid cotton after 60 My of evolution. *Molecular Biology and Evolution* **32**: 1063–1071.

772 **Renny-Byfield S, Kovařík A, Chester M, Nichols RA, Macas J, Novák P, Leitch AR**. **2012**.

773 Independent, rapid and targeted loss of highly repetitive DNA in natural and synthetic

774 allopolyploids of Nicotiana tabacum. *PloS one* **7**: e36963.

775 **Roddy AB, Théroux-Rancourt G, Abbo T, Benedetti JW, Brodersen CR, Castro M, Castro**

776 **S, Gilbride AB, Jensen B, Jiang G-F,** *et al.* **2020**. The scaling of genome size and cell size

777 limits maximum rates of photosynthesis with implications for ecological strategies. *International*

778 *Journal of Plant Sciences* **181**: 75–87.

779 **Sawangproh W, Hedenäs L, Lang AS, Hansson B, Cronberg N**. **2020**. Gene transfer across

780 species boundaries in bryophytes: evidence from major life cycle stages in *Homalothecium*

781 *lutescens* and *H. sericeum*. *Annals of Botany* **125**: 565–579.

782 **Siljak-Yakovlev S, Pustahija F, Solic EM, Bogunic F, Muratovic E, Basic N, Catrice O,**

783 **Brown SC**. **2010**. Towards a genome size and chromosome number database of Balkan flora:

784 C-values in 343 taxa with novel values for 242. *Advanced Science Letters* **3**: 190–213.

785 **Simonin KA, Roddy AB**. **2018**. Genome downsizing, physiological novelty, and the global

786 dominance of flowering plants. *PLOS Biology* **16**: e2003706.

787 **Šímová I, Herben T**. **2012**. Geometrical constraints in the scaling relationships between

788 genome size, cell size and cell cycle length in herbaceous plants. *Proceedings of the Royal*

789 *Society B: Biological Sciences* **279**: 867–875.

790 **Šmarda P, Bureš P**. **2010**. Understanding intraspecific variation in genome size in plants.

791 *Preslia* **82**: 41–61.

792 **Šmarda P, Horová L, Bureš P, Hralová I, Marková M**. **2010**. Stabilizing selection on genome

793 size in a population of *Festuca pallens* under conditions of intensive intraspecific competition.

794 *New Phytologist* **187**: 1195–1204.

795 **Strasburg JL, Rieseberg LH**. **2008**. Molecular demographic history of the annual sunflowers

796 *Helianthus annuus* and *H. petiolaris* - Large effective population sizes and rates of long-term

797 gene flow. *Evolution* **62**: 1936–1950.

798 **Tek AL, Song J, Macas J, Jiang J**. **2005**. Sobo, a recently amplified satellite repeat of potato,

799 and its implications for the origin of tandemly repeated sequences. *Genetics* **170**: 1231–1238.

800 **Twyford AD, Friedman J**. **2015**. Adaptive divergence in the monkey flower *Mimulus guttatus* is

801 maintained by a chromosomal inversion. *Evolution* **69**: 1476–1486.

802 **Vallès J, Canela MÁ, Garcia S, Hidalgo O, Pellicer J, Sánchez-Jiménez I, Siljak-Yakovlev**

803 **S, Vitales D, Garnatje T**. **2013**. Genome size variation and evolution in the family Asteraceae.

804 *Caryologia* **66**: 221–235.

805 **Van't Hof J, Sparrow AH**. **1963**. A relationship between DNA content, nuclear volume, and

806 minimum mitotic cycle time. *Proceedings of the National Academy of Sciences of the United*

807 *States of America* **49**: 897–902.

808 **Vitales D, Álvarez I, Garcia S, Hidalgo O, Nieto Feliner G, Pellicer J, Vallès J, Garnatje T**.

809 **2020**. Genome size variation at constant chromosome number is not correlated with repetitive

810 DNA dynamism in *Anacyclus* (Asteraceae). *Annals of Botany* **125**: 611–623.

811 **Wang X, Gussarova G, Ruhsam M, de Vere N, Metherell C, Hollingsworth PM, Twyford**

812 **AD**. **2018a**. DNA barcoding a taxonomically complex hemiparasitic genus reveals deep

813 divergence between ploidy levels but lack of species-level resolution. *AoB PLANTS* **10**.

814 **Wang W, Mauleon R, Hu Z, Chebotarov D, Tai S, Wu Z, Li M, Zheng T, Fuentes RR, Zhang**

815  **F, *et al.* 2018b**. Genomic variation in 3,010 diverse accessions of Asian cultivated rice. *Nature*

816  **557**: 43–49.


817  **Wang N, McAllister HA, Bartlett PR, Buggs RJA**. **2016**. Molecular phylogeny and genome

818  size evolution of the genus *Betula* (Betulaceae). *Annals of Botany* **117**: 1023–1035.


819  **Weiss-Schneeweiss H, Greilhuber J, Schneeweiss GM**. **2006**. Genome size evolution in

820  holoparasitic *Orobanche* (Orobanchaceae) and related genera. *American Journal of Botany* **93**:

821  148–156.


822  **Wendel JF**. **2015**. The wondrous cycles of polyploidy in plants. *American Journal of Botany*

823  **102**: 1753–1756.


824  **Wong C, Murray BG**. **2012**. Variable changes in genome size associated with different

825  polyploid events in *Plantago* (Plantaginaceae). *Journal of Heredity* **103**: 711–719.


826  **Wulff HD**. **1939**. Chromosomenstudien an der schleswigholsteinischen Angiospermen-Flora.

827  *Berichte der Deutschen Botanischen Gesellschaft* **57**: 84–91.


828  **Zenil-Ferguson R, Ponciano JM, Burleigh JG**. **2016**. Evaluating the role of genome

829  downsizing and size thresholds from genome size distributions in angiosperms. *American*

830  *Journal of Botany* **103**: 1175–1186.


831  **Zhao X-. P, Si Y, Hanson RE, Crane CF, Price HJ, Stelly DM**. **1998**. Dispersed repetitive DNA

832  has spread to new genomes since polyploid formation in cotton. *Genome Research* **8**.

833


834  **SUPPLEMENTAL DATA**

**a**

| 2x | | | | | |
|---|---|---|---|---|---|
| Peak | | Index | Mean | Area | CV% |
| **1** = 2C (*O. sativa*) | | 1.000 | 216.47 | 927 | 2.96 |
| **2** = 2C (*E. anglica*) | | 1.292 | 279.78 | 857 | 2.41 |
| **3** = 4C (*E. anglica*) | | 2.484 | 537.76 | 616 | 2.61 |

**b**

| 4x | | | | | |
|---|---|---|---|---|---|
| Peak | | Index | Mean | Area | CV% |
| **1** = 2C (*O. sativa*) | | 1.000 | 194.71 | 1200 | 2.48 |
| **2** = 2C (*E. nemorosa*) | | 2.281 | 444.22 | 2431 | 1.86 |
| **3** = 4C (*E. nemorosa*) | | 4.498 | 875.77 | 399 | 1.80 |

**c**

**2x - *E. anglica***

| Peak | | Index | Mean | Area | CV% |
|---|---|---|---|---|---|
| **1** = 2C (20151877-B131) | | 1.000 | 223.75 | 540 | 1.73 |
| **2** = 2C (20151862-B126) | | 1.083 | 242.37 | 756 | 2.37 |
| **3** = 4C (20151877-B131) | | 2.094 | 468.50 | 293 | 3.71 |
| **4** = 4C (20151877-B126) | | 2.198 | 491.80 | 345 | 3.61 |

**d**

**4x - *E. arctica***

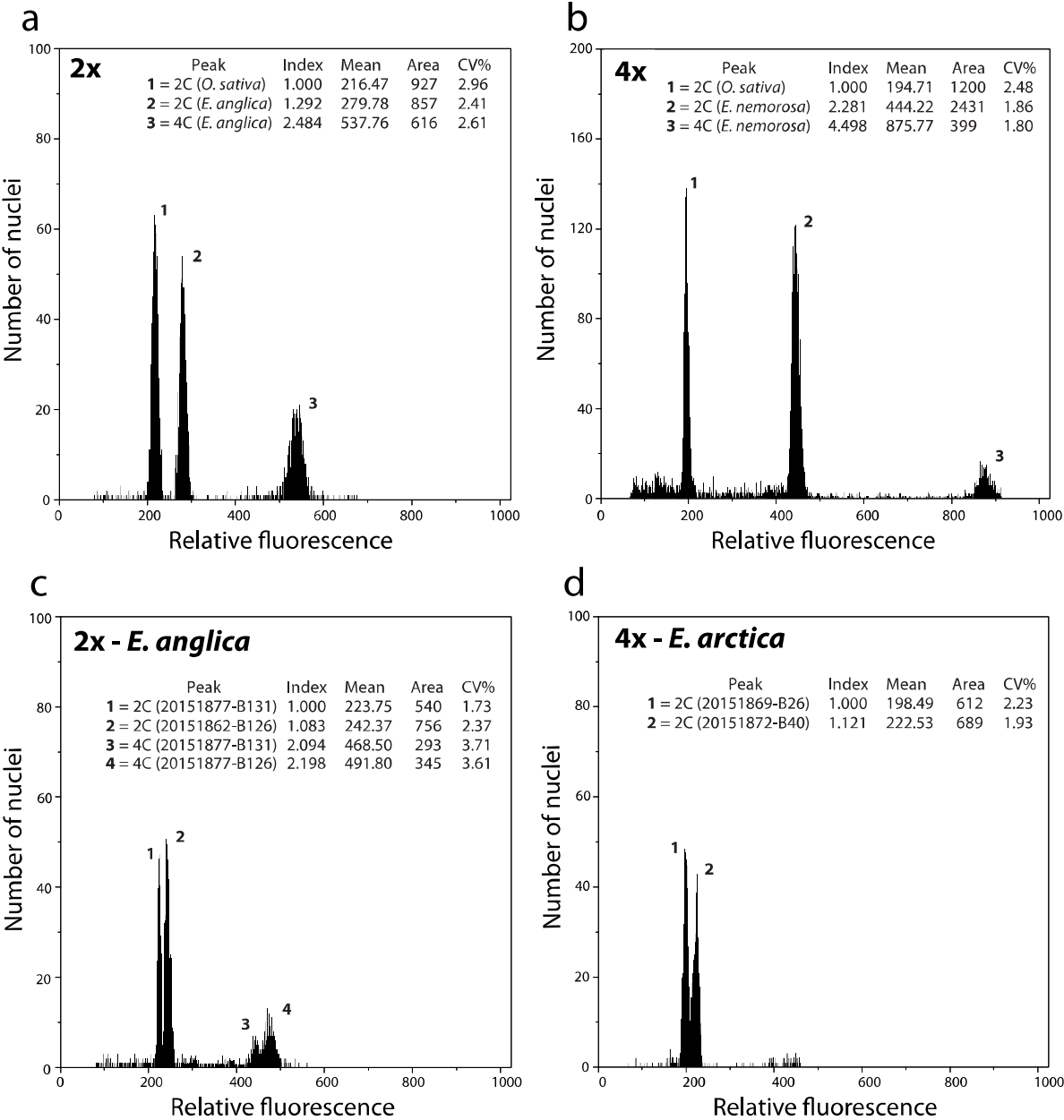| Peak | | Index | Mean | Area | CV% |
|---|---|---|---|---|---|
| **1** = 2C (20151869-B26) | | 1.000 | 198.49 | 612 | 2.23 |
| **2** = 2C (20151872-B40) | | 1.121 | 222.53 | 689 | 1.93 |

835

**Figure S1. Flow cytometry histograms.** A diploid (a) and a tetraploid (b) sample. Intraspecific GS variation in a diploid (c) and a tetraploid (d) species.
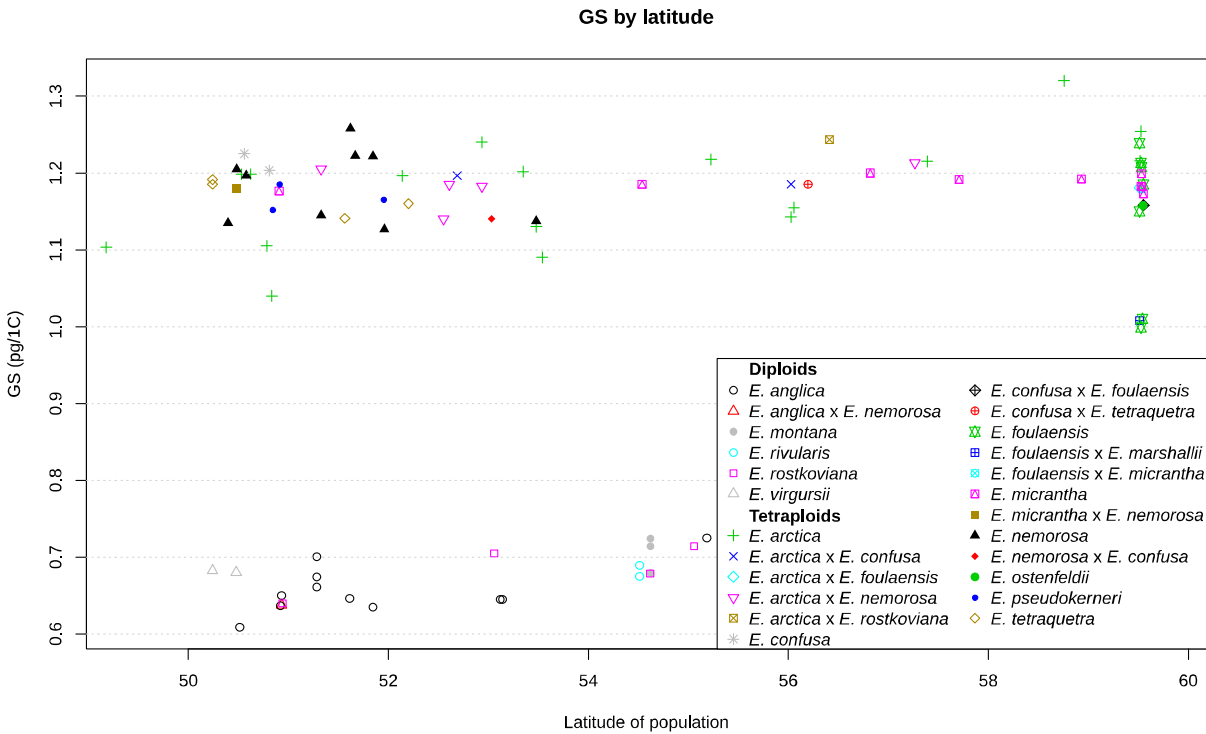
838

33

**GS by latitude**

839

**Figure S2. Genome size plotted against latitude.**

841

**Table S1. Sample information and genome size information**

(Submitted separately)

**Table S2. Details of the whole-genome sequencing data sets generated and genomic proportions of repeat types.**

(Submitted separately)

847