1    # High-throughput genotyping for species identification and diversity assessment

2    # in germplasm collections

3

4    Annaliese S. Mason[a,b], Jing Zhang[c,d], Reece Tollenaere[a,b], Paula Vasquez Teuber[a,b], Jessica Dalton-

5    Morgan[a,b], Liyong Hu[c], Guijun Yan[d,e], David Edwards[a,d], Robert Redden[f], Jacqueline Batley[a,b,d]*

6

7    [a] School of Agriculture and Food Sciences and [b] Centre for Integrative Legume Research, The

8    University of Queensland, Brisbane , 4072, QLD, Australia

9

10   [c] Ministry of Agriculture (MOA) key laboratory of Huazhong Crop Physiology, Ecology and Production,

11   College of Plant Science & Technology, Huazhong Agricultural University, Wuhan, 430070, China.

12

13   [d] School of Plant Biology, Faculty of Science and [e] The UWA Institute of Agriculture, The University of

14   Western Australia, Perth, 6009, WA, Australia

15

16   [f] Australian Grains Genebank, Department of Environment and Primary Industries, Horsham, 3401,

17   VIC , Australia

18

19   * corresponding author; jacqueline.batley@uwa.edu.au; Tel: +61 (0) 7 334 69534; Fax:+ 61 (0)7 336

20   59556

21

24   **Running title:** High throughput germplasm genotyping

**Abstract**

Germplasm collections provide an extremely valuable resource for breeders and researchers. However, misclassification of accessions by species often hinders the effective use of these collections. We propose that use of high-throughput genotyping tools can provide a fast, efficient and cost-effective way of confirming species in germplasm collections, as well as providing valuable genetic diversity data. We genotyped 180 Brassicaceae samples sourced from the Australian Grains Genebank across the recently released Illumina Infinium Brassica 60K SNP array. Of these, 76 were provided on the basis of suspected misclassification and another 104 were sourced independently from the germplasm collection. Presence of the A and C genomes combined with principle components analysis clearly separated *B. rapa*, *B. oleracea*, *B. napus*, *B. carinata* and *B. juncea* samples into distinct species groups. Several lines were further validated using chromosome counts. Overall, 18% of samples (32/180) were misclassified on the basis of species. Within these 180 samples, 23/76 (30%) supplied on the basis of suspected misclassification were misclassified, and 9/105 (9%) of the samples randomly sourced from the Genebank were misclassified. Surprisingly, several individuals were also found to be the product of interspecific hybridisation events. The SNP (Single Nucleotide Polymorphism) array proved effective at confirming species, and provided useful information related to genetic diversity. As similar genomic resources become available for different crops, high-throughput molecular genotyping will offer an efficient and cost-effective method to screen germplasm collections worldwide, facilitating more effective use of these valuable resources by breeders and researchers.

**Introduction**

Natural genetic diversity in crop species is a key resource for agricultural improvement. Genetic variation for cold and heat tolerance, drought and disease resistance as well as other environmental stresses exists in most natural species, but is often lost through domestication and selection for yield and yield-related traits in crops (Day 1973; Hyten *et al.* 2006; Simmonds 1962; Zamir 2001). In order to preserve this useful genetic diversity for later introgression back into crop cultivars and for targeted breeding attempts in crop improvement, "genebanks" and diversity collections exist around the world (Tanksley& McCouch 1997). These collections preserve wild accessions, landraces and cultivars collated from local and international sources, often comprising tens to thousands of lines. Seeds are donated by breeders, collectors and research institutions, and lines are maintained as a resource for future generations.

*Brassica* comprises the largest number of domesticated crop species of any genus, and includes leaf vegetables, oilseeds, condiments and root vegetable crops; such as rapeseed, mustards, cabbage, turnips, broccoli and cauliflower. Numerous species in the wider Brassicaceae can also be hybridised with key crop species within *Brassica*, including the wild radishes (*Raphanus*), woad (*Isatis*) and white mustard (*Sinapis*), as well as the *Brassica* C genome clade of *B. cretica, B. hilarionis*, *B. incana* and *B. macrocarpa*, among others (FitzJohn *et al.* 2007; Harberd& McArthur 1980; Prakash *et al.* 1999; Warwick *et al.* 2003). This potential for hybrid introgression from wild relatives coupled with the extant genetic diversity in the non-cultivated forms of key crop species makes *Brassica* a major feature of genebank collections worldwide. The six cultivated *Brassica* species share an interesting genomic relationship, with three diploids (*B. rapa*, 2n = AA = 20; *B. nigra*, 2n = BB = 16 and *B. oleracea*, 2n = CC = 18) and a set of three allotetraploids each containing two of the three diploid genomes (*B. juncea*, 2n = AABB = 36; *B. napus*, 2n = AACC = 38 and *B. carinata*, 2n = BBCC = 34) (Morinaga 1934; U 1935). Allotetraploid *B. napus* is one of the most agriculturally significant crop

70 species within this genus, with rapeseed and canola contributing to oil production for food and

71 biofuel. However, canola is also the least diverse, with major genetic bottlenecks as a result of only a

72 limited number of hybridisation events between diploid progenitors to form the allotetraploid

73 (Palmer *et al.* 1983), coupled with rigorous selective pressure to achieve "canola-quality" oil for

74 human consumption and enhance yield with the recent emphasis on breeding of oilseeds in this

75 domesticated crop (Cowling 2007). No known wild forms of this species exist (Dixon 2007). Hence, *B.*

76 *napus* in particular is a critical crop species for genetic improvement via introgression of diversity

77 from both wild and domestic diploid relatives, particularly those with which it shares the A and C

78 genomes (*B. rapa*, *B. oleracea*, *B. juncea* and *B. carinata*). Several past breeding attempts have

79 demonstrated the efficacy of this approach in introgressing disease resistance from related species

80 (Navabi *et al.* 2010b; Rygulla *et al.* 2007; Saal *et al.* 2004).

81

82 A major problem with genebank collections is ensuring the accurate identification of species. Many

83 genebanks do not have the resources to assess every line gifted to them for genetic diversity, correct

84 origin and correct species identification. To date, attempts to identify species in germplasm

85 collections have all relied on low-throughput molecular marker genotyping approaches (Dangl *et al.*

86 2001; Ferriol *et al.* 2003; Lee *et al.* 2014; Martin *et al.* 1997; Pradhan *et al.* 2011). However,

87 generation of inexpensive high-throughput molecular marker data is now becoming routine for many

88 genera. We show how the recently released Illumina Infinium Brassica 60K SNP array can be used for

89 rapid species identification in the *Brassica* genus, revealing cases of species misclassification,

90 providing useful genetic diversity information and confirming genome composition in this major

91 agricultural genus.

92

93 **Materials and Methods**

94 *Germplasm*

95   A total of 188 experimental samples (176 lines) were genotyped for this experiment (Supplementary

96   Table 1). A set of 77 samples with suspected species attribution errors and another set of 111

97   independently-obtained samples were sourced from the Australian Grains Genebank (Supplementary

98   Table 1). Forty two additional samples (37 lines) of confirmed species origin were also included in the

99   analysis as controls (Supplementary Table 1). These comprised 22 *B. napus* lines (commercially

100  available canola cultivars from Australia and China), four *B. juncea* lines ("JN9-04", "Purple Leaf

101  Mustard", "Domo" and "Lethbridge"), two *B. carinata* lines ("195923" and "94024", breeding lines in

102  Australia of Ethiopian origin), two *B. oleracea* lines (sequenced accession "TO1000" and commercially

103  available cauliflower "Snowball"), two *B. rapa* lines (sequenced South Korean cultivar "Chiifu" and a

104  commercial "Pak Choy" variety) and five *Raphanus sativus* lines (commercial radish varieties "Cherry

105  Belle", "Long Scarlet", "Mila", "Saxa" and "Scarlet Globe").

106

107  *Genotyping and statistical analyses*

108  DNA was extracted according to methodology detailed in Fulton et al. (1995). All DNA samples were

109  hybridized to an Illumina Infinium Brassica 60K array SNP array released for the *Brassica napus*

110  genome (http://illumina.com; 52157 SNPs) according to manufacturer's instructions. SNP (Single

111  Nucleotide Polymorphism) chips were scanned using an Illumina HiScanSQ and data visualised using

112  Genome Studio V2011.1 (Illumina, Inc., San Diego, CA, USA). A cluster file provided by Agriculture

113  and Agri-Food Canada, Saskatoon, Canada was used to cluster SNPs into genotype groupings (e.g.

114  GG, GT and TT allele calls, which were converted into 0, 1 and 2 scores for subsequent analysis). SNP

115  locations were determined through BLAST comparison with the public *B. rapa* and *B. oleracea*

116  reference genome sequences (Parkin *et al.* 2014; Wang *et al.* 2011); Supplementary Table 2.

117  Percentage SNP calls for each genome were calculated for each sample and this information used to

118  determine the presence or absence of the A and C genomes in the sample.

119

120    Hierarchical clustering and principle components analysis (PCA) were carried out using R version 3.0

121    (The R Project for Statistical Computing). Dendrograms were generated using n = 1000 bootstrap

122    iterations to validate branches, using the "pvclust" function in R package "pvclust". Dendrogram

123    "Height" represents squared Euclidean distance between samples. Missing values were replaced with

124    means for each SNP across the population using R package "gam", function "na.gam.replace". PCA

125    was carried out and output graphs generated using the "dudi.pca" function in R package "ade4".

126

127    *Chromosome counting*

128    Seeds from five experimental lines were germinated on petri dishes under laboratory conditions

129    before harvesting root tip meristems. Root tips were collected and chromosome spreads prepared

130    according to protocols detailed in Mason et al. (2014), using DAPI (4,6-diamidino-2-phenylindole) as a

131    fluorescent stain. Pictures were taken on a Nikon Eclipse E600 microscope with digital camera.

132

133    **Results**

134    *Presence and absence of the A and C genomes*

135    The Illumina Infinium Brassica 60K array comprises 52 157 SNPs. Of these, 10 634 (20.4%) were

136    removed as unreliable or non-specific (consistently amplifying alleles at more than one locus) on the

137    basis of information provided by the Illumina Infinium Brassica 60K cluster file. Of the remaining 41

138    523 SNPs, 44.5% (18 471) were physically located on the *B. rapa* genome (Wang *et al.* 2011) and

139    53.4% (22 155) on the *B. oleracea* genome (Parkin *et al.* 2014). Approximately 12% of these A-

140    genome SNPs also amplified C-genome alleles (in *B. carinata* and *B. oleracea* controls with no A

141    genome), and approximately 23% of these C-genome SNPs also amplified A-genome alleles (in *B.*

142    *rapa* and *B. juncea* controls with no C genome). *Raphanus sativus* samples amplified 13% of alleles on

143    average, with no difference in amplification between the A and C genome SNPs (p = 0.2, Student's t-

144    test).

145

146    A set of 43 control samples (3 *B. rapa*; 6 *B. juncea*, 23 *B. napus,* 2 *B. oleracea*, 4 *B. carinata* and 5

147    *Raphanus sativus*) were run on the Illumina Infinium 60K SNP array. Amplification of A and C genome

148    alleles was assessed in these samples. Clear groups could be distinguished on the basis of A and C

149    genome presence or absence in the controls (Supplementary Figure 1); these groups corresponded

150    to the expected genome presence/absence for each species sample. Of the 188 samples in the

151    experimental population, 59 samples could be classed as "A genome only", 16 samples could be

152    classed as "C genome only", 101 samples could be classed as "A + C genomes" and two samples

153    could be classed as "neither A or C genome present" (Figure 1). An additional seven samples were

154    considered to have failed due to poor quality amplification (removed from further analysis and not

155    included in Figure 1), and another three samples were considered anomalous (included in Figure 1).

156    Two of these samples (R14 and J16) were included in subsequent "A genome only" analyses, and one

157    sample (I2) was discarded from further analysis, leaving 180 samples. On this basis alone, 29/180 of

158    the samples (16%) could be identified to belong to a different species than the one in the genebank

159    records (Supplementary Table 1, Figure 1). Presence of both the A and C genomes also provided a

160    unique identifier for *Brassica napus* samples: 83% of samples (95/115) thought to be *B. napus* were

161    actually *B. napus* (Supplementary Table 1, Figure 1).

162

163    A robust cut-off for sample quality was >75% amplification (an allele call for >75% of SNPs in the A

164    and/or C genome rather than no call reliably indicated genome presence) or <35% amplification in

165    each genome (an allele call for <35% of SNPs in the A and/or C genome reliably indicated genome

166    absence). Samples with 32-57% A and C genome amplification (Supplementary Table 1) also showed

167    random patterns of allele calls and missing data across chromosomes, indicative of unreliable and

168    poor quality SNP data. One of the three samples considered to be anomalous was a putative *B. nigra*

169    sample (I2) that showed 36% A genome and 41% C genome amplification (Figure 1); this may be due

170    to misclassification of this sample coupled with poor quality amplification. The second sample (J16)

171    considered to be anomalous showed 70% A genome amplification and 39% C genome amplification

172    (Figure 1). The third sample considered to be anomalous (putative *B. rapa* sample R14) had 89% A

173    genome presence and 49% C genome presence: on closer inspection of the SNP data, this individual

174    showed presence of some C genome chromosome segments (27 Mbp of C1, all of C2, 7 Mbp of C5,

175    24 Mbp of C6, 30 Mbp of C7 and 39 Mbp of C8). Although material was not available from the

176    individual genotyped, the presence of only 20 chromosomes was confirmed in other individuals from

177    this same line by chromosome counting. Anomalous samples J16 and R14 were retained in our

178    analysis, and sample I2 was discarded.

179

180    *Phylogenetic groupings for species identification*

181    Hierarchical clustering and principle components analysis were performed to separate *B. juncea* and

182    *B. rapa* individuals and *B. carinata* and *B. oleracea* individuals. The *B. juncea* and *B. rapa* group (as

183    deduced from genome presence/absence to have only the A genome) comprised 9 controls and 61

184    experimental individuals. Of the 18 471 SNPs physically mapping to the A genome, 11 983 were

185    polymorphic and amplified in ≥ 90% of the individuals in the population, and were hence used for

186    subsequent analysis. Hierarchical clustering allowed separation of *B. rapa* and *B. juncea* lines, but

187    although species-specific clades were apparent, 100% confidence was not achieved for clade

188    separation (Figure 2; numbers in green and red represent the number of times each branch was in

189    the same position over the 1000 iterations, hence $P<0.05$ = 95 or greater). PCA provided clear

190    separation between *B. rapa* and *B. juncea*, with the first two axes separating two *B. rapa* clades and

191    separating these two groups from *B. juncea* clade, contributing to 18.4% and 13.9% of the variance

192    respectively (Figure 3). Sixty-eight axes were generated, with 48.7% of the variance explained by the

193    first five axes of the PCA.

194

195   The *B. carinata* and *B. oleracea* group as identified by presence of only the C genome consisted of 6

196   control samples and 16 experimental samples. Of the SNP markers mapped to the C genome, 12 794

197   were polymorphic and amplified in ≥ 90% of the individuals in the population, and were hence used

198   for subsequent analysis. Although the *B. carinata* clade fell within the wider *B. oleracea* group, these

199   individuals formed a smaller subgroup with 100% confidence for clade identity using hierarchical

200   clustering (Figure 4). Principle components analysis also showed very clear separation of *B. oleracea*

201   and *B. carinata* samples (first and second axes 41.3% and 13.0% of the variance respectively) and

202   extremely tight grouping of *B. carinata* samples relative to the *B. oleracea* types (Figure 5).

203

204   Overall, 18% of samples (32/180) were misclassified on the basis of species (Table 1). Of the samples

205   suspected to be misclassified, 23/76 (30%) were indeed a species different to the one listed by the

206   Australian Grains Genebank. Of the samples otherwise sourced from the Australian Grains Genebank,

207   9/104 (9%) were misclassified on the basis of species. *B. napus* was observed to be mistaken for each

208   of *B. rapa*, *B. juncea* and *B. carinata*; *B. juncea* was mistaken for *B. rapa* and *B. napus* and *B. rapa* was

209   mistaken for *B. juncea* (Table 1). A complete set of source, species and cultivar/landrace/wild type

210   classifications from the Australian Grains Genebank with confirmed species identifications and SNP

211   genome amplification and heterozygosity results is provided in Supplementary Table 1. Lines were

212   supplied by the Australian Grains Genebank with the label "Advanced cultivar", "Breeder's Line",

213   "Traditional Cultivar/Landrace", "Wild" or "Unknown". Of the 75 samples in the "Advanced cultivar"

214   category, 9 were misclassified (12%). "Traditional cultivar/landraces" had 2/22 samples misclassified

215   (9%) and "Breeder's Line" samples had 2/21 samples misclassified (10%). The single "Wild" sample

216   was also misclassified. "Unknown" samples were misclassified 21 % of the time (11/61).

217

218   *Genetic diversity*

219    Genome diversity within the A genome was assessed in *B. napus*, *B. juncea* and *B. rapa* lines using 13

220    292 polymorphic SNPs amplifying in ≥ 90% of the individuals. Percentage heterozygosity for each

221    individual within the A genome was also calculated using the entire set of A-genome specific SNPs

222    (Supplementary Table 1). C genome diversity was assessed in *B. napus*, *B. oleracea* and *B. carinata*

223    lines using 18 076 SNPs amplifying in ≥ 90% of the individuals and not monomorphic in the

224    population. Percentage heterozygosity for each individual within the C genome was also calculated

225    using the whole set of C-genome specific SNPs (Supplementary Table 1).

226

227    *Brassica rapa* samples putatively from India and Bangladesh based on provenance of samples R05

228    and R21 (Supplementary Table 1, leftmost clade in Figure 3) formed a clearly distinct subgroup when

229    compared to other samples originating from Europe and the rest of Asia. This grouping was not

230    apparent in the first two axes of the PCA of A-genome diversity including the *B. napus* samples

231    (Figure 6). Two outliers were observed on the basis of A-genome diversity using PCA: J06 and J08

232    (Figure 6), which were both reported to be *B. juncea* from China but showed presence of both the A

233    and C genomes; however, using hierarchical clustering analysis these individuals fell within the *B.*

234    *juncea* clade (Supplementary Figure 2). Both individuals had very high A genome heterozygosity (40

235    and 49%) but lower C genome heterozygosity (7 and 21%; Supplementary Table 1).

236

237    As previously observed (Figure 5), the *B. carinata* clade formed a group of tightly-related lines

238    nestled within the *B. oleracea* samples using hierarchical cluster analysis (Supplementary Figure 3).

239    All *B. napus* lines fell outside the *B. oleracea*/*carinata* clade except for three: N019a, N038 and N074

240    (Supplementary Figure 3). Principle components analysis placed N019a within the *B. oleracea*

241    samples, with N038 and N074 in the *B. napus* group but close to *B. oleracea* (Figure 7). N019b, a

242    separately sourced individual of the same accession as N019a, was confirmed to be *B. carinata* due

243    to lack of A genome alleles. N019a had a complete A and C genome, but showed 8.5% heterozygosity

244  in the A genome and 43% C genome heterozygosity, the highest C genome heterozygosity of any

245  experimental *B. napus* sample (Supplementary Table 1). N038 and N074 both had high A- and C-

246  genome heterozygosity (25 – 36% per genome, Supplementary Table 1).

247

248  *Chromosome counting*

249  Chromosome counts were performed for five experimental lines: N067, N089, R05, R14 and R21

250  (Figure 8). Putative *B. napus* sample N067 was confirmed to be *B. juncea* (2n = 36 chromosomes)

251  rather than *B. napus* or *B. rapa*, and putative *B. napus* sample N089 was confirmed to be *B. carinata*

252  (2n = 34 chromosomes) rather than *B. napus* or *B. oleracea*. Each of putative *B. rapa* samples R05,

253  R14 and R21 had 2n = 20 chromosomes, confirming that these plants were *B. rapa*.

254

255  **Discussion**

256  Germplasm collections and genebanks provide an excellent resource for breeders and researchers.

257  However, misclassification of sample genotype and even species is common. Here, we evaluate the

258  use of a high-throughput genotyping technology for the assessment of germplasm collections: the

259  Illumina SNP array, which is increasingly becoming available and cost-effective for many species of

260  interest. We used the Illumina Brassica 60K SNP array for species identification in 180 Brassicaceae

261  samples from the Australian Grains Genebank, a widely used germplasm collection housed in

262  Horsham, Victoria, Australia. The Illumina SNP array provided a quick and effective means to classify

263  species and assess genetic diversity in these samples. A total of 18% of samples were found to be

264  misclassified on the basis of species, and several subpopulations were identified within the various

265  *Brassica* species. A few individuals were also unexpectedly found to result from interspecific

266  hybridisation. This information will prove valuable to future users of this germplasm resource, and

267  validates the use of the Illumina SNP array system for high–throughput genotyping of germplasm

268  collections, particularly in *Brassica*.

269

270 Molecular markers have been used to genotype germplasm collections in the past: SRAP and AFLP

271 markers have been used in cucumber (Ferriol *et al.* 2003), RAPD markers have been used in rice

272 (Martin *et al.* 1997) and SSR markers have been used in grape (Dangl *et al.* 2001) and safflower (Lee

273 *et al.* 2014). High-throughput molecular genotyping is now also starting to be used in major crops: a

274 recent study used genotyping-by-sequencing to characterise lines in the USA national maize inbred

275 seed bank (Romay *et al.* 2013). Problems of species identity within germplasm collections are

276 widespread: in rice, 9/62 (15%) of wild *Oryza* accessions were found to be misclassified; 2/41 grape

277 lines were misclassified, and in another *B. nigra* study using SSR markers, 16/60 (27%) accessions

278 were found to not be *B. nigra* (Pradhan *et al.* 2011). However, older marker technologies are

279 generally not high-throughput, and species identification in germplasm collections using molecular

280 markers has remained out of reach in terms of time and cost until now. In *Brassica* in particular, the

281 high level of homoeology between the A and C genomes, and the presence of multiple species

282 sharing these genomes, can make identification of species-specific alleles difficult (Li *et al.* 2013). In

283 our study, the provision of SNP markers already mapped to the reference genome sequences, a

284 resource which is increasingly available for species of interest, allowed much greater resolution and

285 effectively separated the closely related *Brassica* species.

286

287 We used both Principal Components Analysis and hierarchical clustering to group individuals based

288 on the SNP data results. Importantly, presence of the A genome only, C genome only or both A and C

289 genomes was first used to discriminate *B. napus* samples from *B. juncea*/*B. rapa* and *B. carinata*/*B.*

290 *oleracea*, as *B. napus* samples were not always otherwise 100% distinguishable from *B. juncea*  or *B.*

291 *carinata*. Principal Components Analysis proved more effective at separating species with shared

292 genomes than hierarchical clustering in our analysis. As allopolyploid species *B. carinata*, *B. juncea*

293 and *B. napus* result from a few hybridisation events between diploid progenitor species *B. rapa*, *B.*

294    *nigra* and *B. oleracea* (Arias *et al.* 2014; Kaur *et al.* 2014), the allopolyploid species form less diverse

295    clades nested within the diversity represented by the diploids. To distinguish between *B. juncea* and

296    *B. rapa* and between *B. carinata* and *B. oleracea*, only shared genome information (A or C genome)

297    was available. Hence, hierarchical clustering, which performs pairwise calculations of similarity

298    between samples, may have been less effective at separating species than Principal Components

299    Analysis, which looks at broader correlations and similarities across the data set. Although

300    hierarchical clustering still showed some utility in discriminating between species (Fig. 2, Fig. 4)

301    single-genome Principal Components Analysis is therefore recommended for this purpose in future

302    studies.

303

304    Interestingly, *B. napus* lines in our study were observed to be mistaken for each of *B. carinata*, *B.*

305    *juncea* and *B. rapa*, but only *B. juncea* was commonly mistaken as *B. napus*. However, more *B. napus*

306    lines were used in this experiment than any other species, hence increasing the chance that

307    misclassification errors would be picked up in *B. napus* relative to the other species. Lines sourced as

308    "Traditional cultivar/landraces" or "Unknown" samples may have been expected to be more

309    commonly misclassified than "Advanced Cultivar" or "Breeding Line" samples. However, although

310    "Unknown" samples comprised by far the largest percentage of misclassified samples (11/25), lines

311    sourced as "Advanced Cultivars" were also likely to be misclassified, with a further 9 samples falling

312    into this category. Some of these may have resulted from mislabelling or contamination during seed

313    collection or during seed regeneration of accessions, particularly in the case of commercially

314    available open pollinated (OP) canola cultivars or lines that have passed through many hands before

315    being donated to the Australian Grains Genebank. However, in many cases accurate phenotypic

316    identification of species misclassification was made by the germplasm curators. Samples suspected

317    to be misclassified by the germplasm bank were three times more likely to actually be misclassified

318    on the basis of species (30% as opposed to 9%). In addition, specific recorded notes or remarks

(Supplementary Table 1)  identified the actual species of the sample in a number of instances. For

320 example, N057 was correctly identified as *B. juncea* based on 2010 phenotype data, and likewise

321 N045, N046, N047 and N048 were suspected to be *B. juncea* or *B. rapa* rather than *B. napus* on the

322 basis of phenotype and were confirmed as *B. rapa* by the SNP molecular data. These findings

323 highlight the significance of obtaining phenotypic data wherever possible as a complement to

324 molecular marker results, and support the important role of expert curators in managing germplasm

325 material.

326

327 One of the most surprising and interesting results was the presence in the germplasm collection of

328 several individuals clearly originating from interspecific hybridisation events. Although this is a

329 common method for crop improvement in the *Brassica* genus (Chen *et al.* 2011; Navabi *et al.* 2010b;

330 Rygulla *et al.* 2007; Seyis *et al.* 2003; Zou *et al.* 2011), and all species assessed in this experiment are

331 known to be able to hybridise (FitzJohn *et al.* 2007), lines resulting from interspecific hybridisation

332 events seem unlikely candidates for donation to a germplasm collection, at least without explicit

333 labelling. Hence, it seems likely that these events were spontaneous and originated as a result of

334 cross-contamination during seed bulking processes. We observed one very clear case of interspecific

335 hybridisation in putative *B. rapa* individual R14, which contained a partial C genome in addition to a

336 complete A genome. Confusingly, chromosome counting of another individual resulting from the

337 same seed packet revealed only 20 chromosomes, suggesting either that C genome fragments were

338 still present in a heterozygous state or that only some individuals from this line were carrying these

339 introgressions. Indirect but compelling evidence for hybridisation between *B. juncea* and *B. napus*

340 was obtained for individuals J06 and J08: both were classified as *B. juncea* but also showed presence

341 of a complete C genome; both had much higher A genome heterozygosity than average (45 and 49%)

342 but normal C genome heterozygosity (Supplementary Table 1), and both fell outside the *B. juncea* –

343 *B. napus* groups in the PCA. Individual N019a was also a strong candidate for an interspecific

344  hybridisation event between *B. napus* and *B. carinata*: individual N019b from the same Australian

345  Grains Genebank line but sourced separately was conclusively *B. carinata*, N019a clustered within

346  the *B. oleracea/B. carinata* clade in both the PCA and hierarchical clustering analysis and N019a also

347  had disproportionately high C genome heterozygosity (43%) but normal A genome heterozygosity

348  (9%). Individuals J22 and J05 both also contained an A and a C genome, but grouped strongly with *B.*

349  *carinata* samples using both PCA and hierarchical clustering. These putative interspecific

350  hybridisation events are plausible: accessions in genebanks are often sown in close proximity, and

351  accidental cross-pollination could occur. Hybridisation between the allotetraploid species is relatively

352  easy when carried out by hand pollination (Mason *et al.* 2011) and interspecific hybrids between the

353  allotetraploids are capable of producing seed when self-pollinated (Mason *et al.* 2011) and when

354  back-crossed to the parent species (Chèvre *et al.* 1997; Navabi *et al.* 2010a). Accessions of different

355  *Brassica* species are often grown adjacently during seed regeneration, allowing opportunity for

356  natural cross-pollination to occur.

357

358  High-throughput genotyping using molecular resources such as SNP chip arrays and genotyping-by-

359  sequencing is becoming both readily accessible and cost-effective for large sample sizes and complex

360  crop genomes (Edwards& Batley 2010; Edwards *et al.* 2013). As demonstrated in our study by

361  identification of A- and C-genome-specific SNPs, the availability of reference genome sequences can

362  also dramatically increase the effectiveness of standard molecular marker approaches. We provide

363  validation of the Illumina Infinium Brassica 60K SNP array for species classification in germplasm

364  collections, and suggest that similar high-throughput SNP genotyping approaches should be carried

365  out in future in germplasm collections to support these valuable resources for research and breeding.

366

367  **Acknowledgements**

373

374 **Table 1: Species identity as confirmed by SNP molecular genotyping in a set of *Brassica* samples and related species sourced from the Australian Grains Genebank.**

| Germplasm collection species | Confirmed species | No. samples | % accuracy overall |
|---|---|---|---|
| *B. napus* | *B. napus* | 95 | |
| *B. napus* | *B. rapa* | 8 | |
| *B. napus* | *B. juncea* | 3 | |
| *B. napus* | *B. carinata* | 9 | |
| | **Subtotal** | **115** | **83%** |
| | | | |
| *B. rapa* | *B. rapa* | 20 | |
| *B. rapa* | *B. juncea* | 1 | |
| | **Subtotal** | **21** | **95%** |
| | | | |
| *B. oleracea* | *B. oleracea* | 3 | |
| | **Subtotal** | **3** | **100%** |
| | | | |
| *B. carinata* | *B. carinata* | 3 | |
| | **Subtotal** | **3** | **100%** |
| | | | |
| *B. juncea* | *B. juncea* | 25 | |
| *B. juncea* | *B. rapa* | 2 | |
| *B. juncea* | *B. napus* | 5 | |
| | **Subtotal** | **32** | **77%** |
| | | | |
| *B. nigra* | *B. nigra* | 1 | |
| *B. nigra* | *B. juncea* | 1 | |
| *Sinapis alba* | *Sinapis alba* | 1 | |
| *Sinapis alba* | *B. nigra* | 1 | |
| *Sinapis alba* | *B. carinata* | 1 | |
| *Raphanus sativus* | *B. napus* | 1 | |
| | **Subtotal** | **6** | **33%** |
| | | | |
| | **TOTAL** | **180** | **82%** |

375 **Figure Legends**

376

377 **Figure 1:** Presence of the *Brassica* A and C genomes using SNP markers in a set of Brassicaceae

378 samples sourced from a germplasm collection: 32 putative *B. juncea* samples, 21 putative *B. rapa*

379 samples, 115 putative *B. napus* samples, 3 putative *B. oleracea* samples, 3 putative *B. carinata*

380 samples, 3 putative *B. nigra* samples, 3 putative *Sinapis alba* samples and 1 putative *Raphanus*

381 *sativus* sample. Three anomalous samples are observed outside the tight genome clusters.

382

383 **Figure 2:** Separation of *Brassica rapa* and *B. juncea* samples using A genome SNP data from the

384 Illumina Infinium Brassica 60K array. Dendrogram generated using default hierarchical clustering in

385 package and function "pvclust" in R v 3.0 using n = 1000 iterations; "au" and "bp" refer to the

386 "approximately unbiased" and "bootstrap probability" p-values for each branch. Control samples

387 from confirmed species genotypes are labelled with "Control_" followed by the species and a

388 genotype designation; experimental samples are labelled by a letter representing the supplied

389 species ("J" for *B. juncea*, "R" for *B. rapa*, "I" for *B. nigra*, "N" for *B. napus* (supplied as *B. napus* but

390 containing only an A genome), and "XS" for non-*Brassica, Sinapis alba* (also containing an A

391 genome)). Individual plants from the same genotype are labelled with the same number but different

392 lowercase letters. Chromosome-counted samples are indicated by red stars.

393

394 **Figure 3**: Separation of *B. rapa* and *B. juncea* samples using Principle Components Analysis (first two

395 axes plotted, explaining 18.2% and 13.7% of the variance respectively). Control samples from

396 confirmed species genotypes are labelled with "Control" followed by the species and a genotype

397 designation; experimental samples are labelled by a letter representing the supplied species ("J" for

398 *B. juncea*, "R" for *B. rapa*, "I" for *B. nigra*, "N" for *B. napus* (supplied as *B. napus* but containing only

399 an A genome), and "XS" for non-*Brassica, Sinapis alba* (also containing an A genome)). Individual

400 plants from the same genotype are labelled with the same number but different lowercase letters.

401 Red stars indicate chromosome-counted samples. Individual R014 was anomalous (putatively *B.*

402 *rapa*) with C-genome introgressions in an A-genome background.

403

404 **Figure 4:** Separation of *Brassica oleracea* and *B. carinata* samples using C genome SNP data from the

405 Illumina Infinium Brassica 60K array. Dendrogram generated using default hierarchical clustering in

406 package and function "pvclust" in R v 3.0 using n = 1000 iterations; "au" and "bp" refer to the

407 "approximately unbiased" and "bootstrap probability" p-values for each branch. Control samples

408 from confirmed species genotypes are labelled with "Control_" followed by the species and a

409 genotype designation; experimental samples are labelled by a letter representing the supplied

410 species ("N" for *B. napus* (supplied as *B. napus* but with no A genome), "O" for *B. oleracea*, "C" for *B.*

411 *carinata* and "XS" for non-*Brassica, Sinapis alba*). Individual plants from the same genotype are

412 labelled with the same number but different lowercase letters. A chromosome-counted sample is

413 indicated with a red star.

414

415 **Figure 5:** Separation of *B. oleracea* and *B. carinata* samples using Principle Components Analysis (first

416 two axes plotted, explaining 41.3% and 13.0% of the variance respectively). Control samples from

417 confirmed species genotypes are labelled with "Control" followed by the species and a genotype

418 designation; experimental samples are labelled by a letter representing the supplied species ("N" for

419 *B. napus* (supplied as *B. napus* but with no A genome), "O" for *B. oleracea*, "C" for *B. carinata* and

420 "XS" for non-*Brassica, Sinapis alba*). Individual plants from the same genotype are labelled with the

421 same number but different lowercase letters. The red star indicates a chromosome-counted sample.

422

423    **Figure 6:** A genome diversity as assessed by Principle Components Analysis of Illumina Infinium 60k

424    Brassica array data in a set of 31 A-genome controls of known species origin and 162 *B. rapa*, *B.*

425    *juncea* and *B. napus* samples found to contain an A genome and originating from the Australian

426    Grains Genebank. Experimental samples are labelled by a letter representing the supplied species

427    ("J" for *B. juncea* and "R" for *B. rapa*).

428

429    **Figure 7:** C genome diversity as assessed by principle components analysis of Illumina Infinium 60k

430    Brassica array data from a set of 29 C-genome controls of known species origin (2 *B. oleracea*, 4 *B.*

431    *carinata* and 23 *B. napus*) and 117 *B. carinata*, *B. oleracea* and *B. napus* samples all containing a C

432    genome and originating from the Australian Grains Genebank. Control samples from confirmed

433    species genotypes are labelled with "Control" followed by the species and a genotype designation;

434    experimental samples are labelled by a letter representing the supplied species: "N" for *B. napus*, "O"

435    for *B. oleracea* and "J" for *B. juncea* (supplied as *B. juncea* but containing an A and a C genome and

436    hence actually *B. napus*). Individual plants from the same genotype are labelled with the same

437    number but different lowercase letters.

438

439    **Figure 8:** Chromosome counts for two putative *Brassica napus* plants (N089 and N067) showing 2n =

440    34 (*B. carinata*) and 2n = 36 (*B. juncea*) respectively; and three *B. rapa* individuals (R05, R14 and R21)

441    showing 2n = 20. Bar = 10 µm

442

443    **References**

444    Arias T, Beilstein MA, Tang M, McKain MR, Pires JC (2014) Diversification times among *Brassica*

445        (Brassicaceae) crops suggest hybrid formation after 20 million years of divergence. *American*

446        *Journal of Botany* **101**, 86-91.

447      Chen S, Nelson MN, Chèvre A-M*, et al.* (2011) Trigenomic bridges for *Brassica* improvement. *Critical*

448          *Reviews in Plant Sciences* **30**, 524-547.

449      Chèvre AM, Barret P, Eber F*, et al.* (1997) Selection of stable *Brassica napus-B. juncea* recombinant

450          lines resistant to blackleg (*Leptosphaeria maculans*). 1. Identification of molecular markers,

451          chromosomal and genomic origin of the introgression. *Theoretical and Applied Genetics* **95**,

452          1104-1111.

453      Cowling WA (2007) Genetic diversity in Australian canola and implications for crop breeding for

454          changing future environments. *Field Crops Research* **104**, 103-111.

455      Dangl GS, Mendum ML, Prins BH*, et al.* (2001) Simple sequence repeat analysis of a clonally

456          propagated species: A tool for managing a grape germplasm collection. *Genome* **44**, 432-438.

457      Day PR (1973) Genetic variability of crops. *Annual Review of Phytopathology* **11**, 293-312.

458      Dixon GR (2007) Vegetable Brassicas and related crucifers. In: *Crop production science in horticulture*

459          *series* (eds. Atherton J, Rees H). CAB International, Oxfordshire, UK.

460      Edwards D, Batley J (2010) Plant genome sequencing: applications for crop improvement. *Plant*

461          *Biotechnology Journal* **8**, 2-9.

462      Edwards D, Batley J, Snowdon RJ (2013) Accessing complex crop genomes with next-generation

463          sequencing. *Theoretical and Applied Genetics* **126**, 1-11.

464      Ferriol M, Pico B, Nuez F (2003) Genetic diversity of a germplasm collection of *Cucurbita pepo* using

465          SRAP and AFLP markers. *Theoretical and Applied Genetics* **107**, 271-282.

466      FitzJohn RG, Armstrong TT, Newstrom-Lloyd LE, Wilton AD, Cochrane M (2007) Hybridisation within

467          *Brassica* and allied genera: evaluation of potential for transgene escape. *Euphytica* **158**, 209-

468          230.

469      Fulton TM, Chunwongse J, Tanksley SD (1995) Microprep protocol for extraction of DNA from tomato

470          and other herbaceous plants. *Plant Molecular Biology Reporter* **13**, 207-209.

471    Harberd DJ, McArthur ED (1980) Meiotic analysis of some species and genus hybrids in the

472        Brassiceae. In: *Brassica Crops and Wild Allies: Biology and Breeding* (ed. Tsunoda S, Hinata,

473        K., Gomez-Campo, C.), pp. 65-87. Japan Scientific Societies Press, Tokyo.

474    Hyten DL, Song QJ, Zhu YL*, et al.* (2006) Impacts of genetic bottlenecks on soybean genome diversity.

475        *Proceedings of the National Academy of Sciences of the United States of America* **103**, 16666-

476        16671.

477    Kaur P, Banga S, Kumar N*, et al.* (2014) Polyphyletic origin of *Brassica juncea* with *B. rapa* and *B. nigra*

478        (Brassicaceae) participating as cytoplasm donor parents in independent hybridization events.

479        *American Journal of Botany* **101**, 1157-1166.

480    Lee GA, Sung JS, Lee SY*, et al.* (2014) Genetic assessment of safflower (*Carthamus tinctorius* L.)

481        collection with microsatellite markers acquired via pyrosequencing method. *Molecular*

482        *Ecology Resources* **14**, 69-78.

483    Li HT, Younas M, Wang XF*, et al.* (2013) Development of a core set of single-locus SSR markers for

484        allotetraploid rapeseed (*Brassica napus* L.). *Theoretical and Applied Genetics* **126**, 937-947.

485    Martin C, Juliano A, Newbury HJ*, et al.* (1997) The use of RAPD markers to facilitate the identification

486        of *Oryza* species within a germplasm collection. *Genetic Resources and Crop Evolution* **44**,

487        175-183.

488    Mason AS, Nelson MN, Takahira J*, et al.* (2014) The fate of chromosomes and alleles in an

489        allohexaploid *Brassica* population. *Genetics* **197**, 273-283.

490    Mason AS, Nelson MN, Yan GJ, Cowling WA (2011) Production of viable male unreduced gametes in

491        *Brassica* interspecific hybrids is genotype specific and stimulated by cold temperatures. *BMC*

492        *Plant Biology* **11**, 103.

493    Morinaga T (1934) Interspecific hybridisation in *Brassica* VI. The cytology of $F_1$ hybrids of *B. juncea*

494        and *B. nigra*. *Cytologia* **6**, 62-67.

495  Navabi ZK, Parkin IA, Pires JC*, et al.* (2010a) Introgression of B-genome chromosomes in a doubled

496        haploid population of *Brassica napus* x *B. carinata*. *Genome* **53**, 619-629.

497  Navabi ZK, Strelkov SE, Good AG, Thiagarajah MR, Rahman MH (2010b) *Brassica* B-genome resistance

498        to stem rot (*Sclerotinia sclerotiorum*) in a doubled haploid population of *Brassica napus* x

499        *Brassica carinata*. *Canadian Journal of Plant Pathology-Revue Canadienne De*

500        *Phytopathologie* **32**, 237-246.

501  Palmer JD, Shields CR, Cohen DB, Orten TJ (1983) Chloroplast DNA evolution and the origin of

502        amphidiploid *Brassica* species. *Theoretical and Applied Genetics* **65**, 181-189.

503  Parkin IA, Koh C, Tang H*, et al.* (2014) Transcriptome and methylome profiling reveals relics of

504        genome dominance in the mesopolyploid *Brassica oleracea Genome Biology* **15**, R77.

505  Pradhan A, Nelson MN, Plummer JA, Cowling WA, Yan GJ (2011) Characterization of *Brassica nigra*

506        collections using simple sequence repeat markers reveals distinct groups associated with

507        geographical location, and frequent mislabelling of species identity. *Genome* **54**, 50-63.

508  Prakash S, Takahata Y, Kirti PB, Chopra VL (1999) Cytogenetics. In: *Biology of Brassica coenospecies*

509        (ed. Gómez-Campo C), pp. 59-106. Elsevier Science B.V., Amsterdam.

510  Romay MC, Millard MJ, Glaubitz JC*, et al.* (2013) Comprehensive genotyping of the USA national

511        maize inbred seed bank. *Genome Biology* **14**.

512  Rygulla W, Friedt W, Seyis F*, et al.* (2007) Combination of resistance to *Verticillium longisporum* from

513        zero erucic acid *Brassica oleracea* and oilseed *Brassica rapa* genotypes in resynthesized

514        rapeseed (*Brassica napus*) lines. *Plant Breeding* **126**, 596-602.

515  Saal B, Brun H, Glais I, Struss D (2004) Identification of a *Brassica juncea*-derived recessive gene

516        conferring resistance to *Leptosphaeria maculans* in oilseed rape. *Plant Breeding* **123**, 505-

517        511.

518    Seyis F, Snowdon RJ, Luhs W, Friedt W (2003) Molecular characterization of novel resynthesized

519          rapeseed (*Brassica napus*) lines and analysis of their genetic diversity in comparison with

520          spring rapeseed cultivars. *Plant Breeding* **122**, 473-478.

521    Simmonds NW (1962) Variability in crop plants, its use and conservation. *Biological Reviews of the*

522          *Cambridge Philosophical Society* **37**, 422-&.

523    Tanksley SD, McCouch SR (1997) Seed banks and molecular maps: Unlocking genetic potential from

524          the wild. *Science* **277**, 1063-1066.

525    U N (1935) Genome-analysis in *Brassica* with special reference to the experimental formation of *B.*

526          *napus* and peculiar mode of fertilization. *Japanese Journal of Botany* **7**, 389-452.

527    Wang XW, Wang HZ, Wang J*, et al.* (2011) The genome of the mesopolyploid crop species *Brassica*

528          *rapa*. *Nature Genetics* **43**, 1035-1039.

529    Warwick SI, Simard M-J, Légère A*, et al.* (2003) Hybridization between transgenic *Brassica napus* L.

530          and its wild relatives: *Brassica rapa* L., *Raphanus raphanistrum* L., *Sinapis arvensis* L., and

531          *Erucastrum gallicum* (Willd.) O.E. Schulz. *Theoretical and Applied Genetics* **107**, 528-539.

532    Zamir D (2001) Improving plant breeding with exotic genetic libraries. *Nature Reviews Genetics* **2**,

533          983-989.

534    Zou J, Fu DH, Gong HH*, et al.* (2011) *De novo* genetic variation associated with retrotransposon

535          activation, genomic rearrangements and trait variation in a recombinant inbred line

536          population of *Brassica napus* derived from interspecific hybridization with *Brassica rapa*.

537          *Plant Journal* **68**, 212-224.

538

539

540 **Data Accessibility**

541 The Illumina Infinium Brassica 60K SNP array used in this analysis can be obtained from Illumina Inc.

542 (http://www.illumina.com/). Summary information for each Australian Germplasm Genebank

543 accession used in this analysis is provided in Supplementary Table 1. Genotype data and SNP

544 information is provided in Supplementary Table 2 and this data is also available *via* the Dryad data

545 repository (doi:10.5061/dryad.c3g5r). Seeds for each of the lines used can be obtained from the

546 Australian Germplasm Genebank. PCA and hierarchical clustering analyses were performed using the

547 R base software and packages "pvclust", "ade4" and "gam" freely available from the R Project for

548 Statistical Computing (http://www.r-project.org/).

549

550 **Author Contributions**

551 JB, DE and BR conceptualised the study. JB managed the project. BR, GY, JZ and LH contributed

552 material. RT and JZ grew up seeds and extracted DNA. PVT carried out chromosome counting. JDM

553 ran the SNP chip. ASM analysed the SNP chip data, generated the figures and tables and wrote the

554 paper. JB, DE, BR and GY critically revised the manuscript. All authors have read and approved the

555 final version of the manuscript.

556

557

558   **Supporting Information**

559

560   **Supplementary Figure 1:** Presence of the *Brassica* A and C genomes in a set of known control

561   samples using SNP markers: 3 *B. rapa* (A genome only); 6 *B. juncea* (A genome only), 23 *B. napus*

562   (A+C genomes)*,* 2 *B. oleracea* (C genome only), 4 *B. carinata* (C genome only) and 5 *Raphanus sativus*

563   (neither genome)*.*

564

565   **Supplementary Figure 2:** A genome diversity as assessed by hierarchical clustering of Illumina

566   Infinium 60k Brassica array data in a set of 31 controls of known species origin and 162 *B. rapa*, *B.*

567   *juncea* and *B. napus* lines originating from the Australian Grains Genebank. Control samples from

568   confirmed species genotypes are labelled with "Control" followed by the species and a genotype

569   designation; experimental samples are labelled by a letter representing the supplied species ("J" for

570   *B. juncea*, "R" for *B. rapa*, "I" for *B. nigra*, "N" for *B. napus* and "XS" for non-*Brassica, Sinapis alba*

571   (also containing an A genome)). Individual plants from the same genotype are labelled with the same

572   number but different lowercase letters. Red stars indicate chromosome-counted samples.

573   Chromosome-counted lines are indicated by red stars, and samples of interest are indicated using

574   blue four-pointed stars.

575

576   **Supplementary Figure 3:** C genome diversity as assessed by hierarchical clustering of Illumina

577   Infinium 60k Brassica array data from a set of 29 controls of known species origin and 117 *B.*

578   *carinata*, *B. oleracea* and *B. napus* lines originating from the Australian Grains Genebank. Control

579   samples from confirmed species genotypes are labelled with "Control" followed by the species and a

580   genotype designation; experimental samples are labelled by a letter representing the supplied

581   species ("N" for *B. napus*, "O" for *B. oleracea*, "C" for *B. carinata*, "J" for *B. juncea* (but containing

582   both an A and C genome and hence actually *B. napus*) and "XS" for non-*Brassica, Sinapis alba*).

583    Individual plants from the same genotype are labelled with the same number but different lowercase

584    letters. Chromosome-counted lines are indicated by red stars, and samples of interest are indicated

585    using blue four-pointed stars.

586

587    **Supplementary Table 1:** Information for the set of 188 experimental samples sourced from the

588    Australian Grains Genebank: sample identification numbers, provided information, genome

589    amplification results and species re-classifications based on SNP analyses.

590

591    **Supplementary Table 2:** SNP molecular genotyping data and information.

592