

Feral populations of *Brassica oleracea* along Atlantic coasts in western Europe

Journal:	<i>Ecology and Evolution</i>
Manuscript ID	ECE-2020-04-00487.R1
Wiley - Manuscript type:	Original Research
Date Submitted by the Author:	23-Jul-2020
Complete List of Authors:	Mittelll, Elizabeth; University of Glasgow, IBAHCM; University of Saint Andrews, Centre for Biodiversity Cobbold, Christina; University of Glasgow, School of Mathematics and Statistics Ijaz, Umer; University of Glasgow, School of Engineering Kilbride, Elizabeth; University of Glasgow, IBAHCM Moore, Karen; University of Exeter, University of Exeter Sequencing Service Mable, Barbara; University of Glasgow, Institute of Biodiversity, Animal Health and Comparative Medicine
Category:	Evolutionary Ecology
Habitat:	Terrestrial
Organism:	Plants
Approach:	Molecular Genetics
Abstract:	<p>There has been growing emphasis on the role that crop wild relatives might play in supporting highly selected agriculturally valuable species in the face of climate change. In species that were domesticated many thousands of years ago, distinguishing wild populations from escaped feral forms can be challenging, but reintroducing variation from either source could supplement current cultivated forms. For economically important cabbages (Brassicaceae: <i>Brassica oleracea</i>), "wild" populations occur throughout Europe but little is known about their genetic variation or potential as resources for breeding more resilient crop varieties. The main aim of this study was to characterise the population structure of geographically isolated wild cabbage populations along the coasts of the UK and Spain, including the Atlantic range edges. Double-digest restriction-site associated DNA sequencing was used to sample individual cabbage genomes, assess the similarity of plants from 20 populations, and explore environment-genotype associations across varying climatic conditions. Interestingly, there were no indications of isolation-by-distance; several geographically close populations were genetically more distinct from each other than to distant populations. Furthermore, several distant populations shared genetic ancestry, which could indicate that they were established by escapees of similar source cultivars. However, there were signals of local adaptation to different environments, including a possible relationship between genetic diversity</p>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

	and soil pH. Overall, these results highlight wild cabbages in the Atlantic region as an important genetic resource worthy of further research into their relationship with existing crop varieties.
Note: The following files were submitted by the author for peer review, but cannot be converted to PDF. You must view these files (e.g. movies) online.	
dataset_1_sumstats.tsv dataset_2_sumstats.tsv dataset_3_sumstats.tsv	



1
2
3
4 **Title: Feral populations of *Brassica oleracea* along Atlantic coasts in western Europe**
5
6
7

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

2
3 Running title: Feral *Brassica oleracea* in western Europe
4

5 Elizabeth A. Mittell^{1,2,*}, Christina A. Cobbold^{3,4}, Umer Zeeshan Ijaz⁵, Elizabeth A. Kilbride¹, Karen A.
6 Moore⁶ & Barbara K. Mable^{1,4,†}

7
8 ¹Institute of Biodiversity, Animal Health and Comparative Medicine, University of Glasgow, UK; ²School
9 of Biology, University of St Andrews, UK; ³School of Mathematics and Statistics, University of Glasgow,
10 UK; ⁴The Boyd Orr Centre for Population and Ecosystem Health, University of Glasgow, Glasgow, UK;
11 ⁵School of Engineering, University of Glasgow, UK; ⁶Exeter Sequencing Service, University of Exeter,
12 UK

13
14 Corresponding authors: Elizabeth A. Mittell, *em294@st-andrews.ac.uk, e.mittell@gmail.com & Barbara
15 K. Mable, † Barbara.Mable@glasgow.ac.uk
16

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract

There has been growing emphasis on the role that crop wild relatives might play in supporting highly selected agriculturally valuable species in the face of climate change. In species that were domesticated many thousands of years ago, distinguishing wild populations from escaped feral forms can be challenging, but reintroducing variation from either source could supplement current cultivated forms. For economically important cabbages (Brassicaceae: *Brassica oleracea*), “wild” populations occur throughout Europe but little is known about their genetic variation or potential as resources for breeding more resilient crop varieties. The main aim of this study was to characterise the population structure of geographically isolated wild cabbage populations along the coasts of the UK and Spain, including the Atlantic range edges. Double-digest restriction-site associated DNA sequencing was used to sample individual cabbage genomes, assess the similarity of plants from 20 populations, and explore environment-genotype associations across varying climatic conditions. Interestingly, there were no indications of isolation-by-distance; several geographically close populations were genetically more

1
2
3
4 30 distinct from each other than to distant populations. Furthermore, several distant populations shared
5
6 31 genetic ancestry, which could indicate that they were established by escapees of similar source culti-
7
8 32 vars. However, there were signals of local adaptation to different environments, including a possible
9
10 33 relationship between genetic diversity and soil pH. Overall, these results highlight wild cabbages in the
11
12 34 Atlantic region as an important genetic resource worthy of further research into their relationship with
13
14 35 existing crop varieties.

15
16 36 **Keywords:** *Brassica oleracea*, feral populations, crop wild relatives, isolation-by-distance, environment-
17
18 37 genotype associations, domestication

21 22 38 **Introduction**

23
24
25 39 Domestication was an important transition within human societies, which allowed the rise of civilisations
26
27 40 (Diamond, 2002). Whilst vital for human success, there have been evolutionary consequences for the
28
29 41 domesticated organisms. In crop plants, the selection of 'domestication traits' has led to many desired
30
31 42 changes in physiological, morphological and life-history traits compared to their wild relatives (Milla,
32
33 43 Osborne, Turcotte, & Violle, 2015; Purugganan & Fuller, 2009). However, traits that are correlated with
34
35 44 those selected for (directly or indirectly) can also influence phenotypes via pleiotropic effects (Conner,
36
37 45 2002) and linkage disequilibrium (Falconer & Mackay, 1996). These genetic constraints and narrow
38
39 46 population bottlenecks can have unintended genetic consequences for crop plants, particularly elite
40
41 47 lines that are the result of intense artificial selection; e.g., reduced genetic diversity, increased genetic
42
43 48 drift and increased deleterious allele frequencies (Rauf, Teixeira da Silva, Khan, & Naveed, 2010; von
44
45 49 Wettberg et al., 2018). It is also likely that crop lines are constrained to some extent by the environment
46
47 50 within which they were originally domesticated. Therefore, to continue to utilise crop plants successfully,
48
49 51 it is important to understand both the genetic consequences of domestication, and where it occurred.

50
51 52 A classic example of domestication can be found in the commercially valuable species, *Brassica ol-*
52
53 53 *eracea* (recognised by Darwin, 1859; Walley et al., 2012). This single species contains a huge amount
54
55 54 of morphological diversity in cultivated varieties that has been around since at least the 1st Century (e.g.,
56
57 55 kale, kohlrabi, broccoli, Brussels sprouts and cauliflower; Maggioni, von Bothmer, Poulsen, & Lipman,
58
59 56 2018); the same morphological extremes are not found in wild populations. The origin of domesticated
60

1
2
3
4 57 *B. oleracea* crops and the 'wild' or 'feral' status of populations, found throughout the UK and along the
5
6 58 Atlantic coasts of north-western Europe (Raybould, Mogg, Clarke, Gliddon, & Gray, 1999), has been de-
7
8 59 bated in the literature (Allender, Allainguillaume, Lynn, & King, 2007; Gómez-Campo & Prakash, 1999;
9
10 60 Maggioni, 2015; Mitchell, 1976). Initially it was thought that different cultivars were independently do-
11
12 61 mesticated from wild populations on European Atlantic coasts (e.g., Spanish cabbage varieties were
13
14 62 domesticated from local wild Spanish populations; Gómez-Campo & Prakash, 1999), and that early
15
16 63 domesticates were introduced to and diversified within the Mediterranean region around 3,000 to 4,000
17
18 64 years ago (Allender et al., 2007). Information was limited when this hypothesis was favoured (Allender et
19
20 65 al., 2007; Gómez-Campo & Prakash, 1999), although there was already conflicting evidence (Mitchell,
21
22 66 1976). For example, Mitchell (1976) found that the locations of ancient human settlements and modern
23
24 67 *B. oleracea* populations coincided along UK coasts, providing a potential source of escapees from do-
25
26 68 mestic settings (agriculture or gardens) that could have established feral populations. This alternative
27
28 69 hypothesis that *B. oleracea* originated elsewhere and escaped into the wild in the Atlantic region has
29
30 70 been supported by recent linguistic and historical research (Maggioni, 2015; Maggioni et al., 2018).
31
32 71 Maggioni (2015) suggested that the most plausible hypothesis is that *B. oleracea* was domesticated in
33
34 72 the Mediterranean region, before being moved across Europe by people, where escaped plants estab-
35
36 73 lished now naturalised populations. However, the genetic status of *B. oleracea* in the Atlantic region is
37
38 74 still an open question (*B. oleracea* is classified as a native species in the UK and an alien species in
39
40 75 Spain; *Euro+Med PlantBase*, 2020).

41
42 76 The ease with which cultivated and wild *B. oleracea* plants can introgress is an issue for interpreting
43
44 77 variation within the *B. oleracea* species complex, as past hybridisation can obscure phylogeographic sig-
45
46 78 nals (Allender et al., 2007). However, for crop breeding purposes a close genetic relationship between
47
48 79 wild populations and domesticated cultivars may be seen as an advantage; higher genetic similarity
49
50 80 could make it easier to introgress adaptive traits from the wild into cultivated varieties (Hoisington et
51
52 81 al., 1999). An alternative view is that if these populations are feral they would have experienced the
53
54 82 same domestication bottleneck as many cultivars (von Wettberg et al., 2018), and therefore they may
55
56 83 not be the important reservoirs of genetic diversity that crop wild relatives are typically assumed to be.
57
58 84 Compared to domestication, feralization is under-investigated; however modern genomic data are al-

1
2
3
4 85 lowing its occurrence to be identified and consequences better understood (see examples in Henriksen,
5
6 86 Gering, & Wright, 2018). Despite the agricultural importance of *B. oleracea*, there has not yet been a
7
8 87 comprehensive genetic analysis of wild populations in the Atlantic region that would allow assessment
9
10 88 of their utility as sources of variation for cultivation.

11
12 89 Escaped plants can be thought of as ‘invasive’ species, which are defined as those that became
13
14 90 established after introduction outside of the biogeographic region within which they evolved (Prentis,
15
16 91 Wilson, Dormontt, Richardson, & Lowe, 2008). However, it is not always clear where these ‘native’
17
18 92 regions are located, as is the case of *B. oleracea*, or why certain species are successful where others
19
20 93 are not. Furthermore, wild populations of *B. oleracea* do not have the characteristics that are thought to
21
22 94 be important for successful establishment in novel locations (i.e. ‘invasive traits’; Funk, Standish, Stock,
23
24 95 & Valladares, 2016). For example, wild *B. oleracea* are: perennials rather than annuals, woody rather
25
26 96 than herbaceous, relatively slow- rather than fast-growing, and predominantly outcrossing rather than
27
28 97 selfing. Self-fertilisation in plants is inhibited by polymorphic self-incompatibility (SI) recognition systems
29
30 98 where haplotype blocks encode distinct proteins for pollen-pistil recognition (Charlesworth, Vekemans,
31
32 99 Castric, & Glémin, 2005). A strong SI system exists in *B. oleracea* (a single-locus system with over
33
34 100 60 alleles; Raybould et al., 1999), making them predominantly self-incompatible (Kitashiba & Nasrallah,
35
36 101 2014; Walley et al., 2012; Yousef, Mueller, Börner, & Schmid, 2018). Development of self-compatible
37
38 102 lines can aid in propagation of cultivated forms (e.g., Xiao et al., 2019), but reduce adaptive potential
39
40 103 to changing environmental conditions. Therefore, even if the “wild” populations include escaped forms,
41
42 104 retention of a wide range of self-incompatibility alleles could be used to enhance the potential of breeding
43
44 105 strategies designed to maintain heterosis.

45
46 106 Currently too little is known about levels of genetic variation and population structure in wild *B. ol-*
47
48 107 *eracea* populations to fully assess the potential for use of plants from different regions to supplement
49
50 108 crop diversity. Population structure and within population genetic diversity are impacted by gene flow,
51
52 109 which occurs via pollen and seeds in plants (Scheepens, Frei, Armbruster, & Stöcklin, 2012; Slatkin,
53
54 110 1987). The main pollinators of *B. oleracea* are bees that fly short distances between plants (average
55
56 111 2 m; Raybould et al., 1999). Seed dispersal was previously thought to be limited to approximately 4 m
57
58 112 (Watson-Jones, Maxted, & Ford-Lloyd, 2006). However, Wichmann et al. (2009) found that wind can
59
60

1
2
3
4 113 spread seeds up to 250 m, and that rare-long distance dispersal events of up to 10 km could occur
5
6 114 if seeds became attached to people's shoes. Therefore, although gene flow may be limited between
7
8 115 geographically close populations leading to high genetic structuring in some instances, in other cases,
9
10 116 such as where plants grow close (0 - 4 m) to well used coastal paths, gene flow might be greater than
11
12 117 expected. Genetic diversity estimates have been made in some *B. oleracea* populations within the
13
14 118 Atlantic region (e.g., Table 1), but the northern edge (Scotland) has not been investigated. A correla-
15
16 119 tion between genetic distance and geographic distance in wild *B. oleracea* populations was found in
17
18 120 some studies (Raybould et al., 1999; Sánchez-Yélamo, 2014) but not others (Christensen et al., 2011;
19
20 121 Watson-Jones et al., 2006). Interestingly, Watson-Jones et al. (2006) also considered some environ-
21
22 122 mental variables and found that higher soil pH was associated with lower genetic diversity in English
23
24 123 and Welsh populations. The inconsistency in previous studies could be due to the varying spatial scales
25
26 124 and molecular markers used. However, overall, these results highlight the uncertainty in the status and
27
28 125 genetic contents of wild *B. oleracea* populations in the Atlantic region, as well as the potential effect of
29
30 126 environment on the plant genetics. Filling these knowledge gaps could provide important insights into
31
32 127 these crop wild relatives for agricultural use.

33
34 128 *Brassica oleracea* is a good model for investigating the genetic resources available (e.g., the extent
35
36 129 of genetic diversity and local adaptation) in a potentially feral crop wild relative because it is diploid
37
38 130 and a reference genome is available (Liu et al., 2014). Therefore, compared to other crop species
39
40 131 (e.g., polyploids) genetic analyses are simpler. For many questions whole-genome sequencing is un-
41
42 132 necessary (Rockman, 2012) and reduced-representation methods, such as double-digest restriction
43
44 133 associated DNA sequencing (ddRADseq), are sufficient to: assess genetic diversity within and between
45
46 134 populations (Andrews, Good, Miller, Luikart, & Hohenlohe, 2016); determine population genetic struc-
47
48 135 turing (Gao et al., 2017); and investigate potential associations between genotypes and environmental
49
50 136 variables (Forester, Lasky, Wagner, & Urban, 2018). Therefore, ddRADseq is an appropriate method
51
52 137 for considering the genetic resources in, and local adaptation of, *B. oleracea* populations across their
53
54 138 Atlantic range.

55
56 139 Overall, current knowledge on genetic variation of *B. oleracea* in wild populations is patchy in geo-
57
58 140 graphic coverage and based on outdated molecular genetic techniques (Table 1). Therefore, this study
59
60

1
2
3
4
5 141 combined modern genetic techniques and the reference genome available for this species to increase
6
7 142 the power to detect differences among populations across a broad geographic range. The following
8
9 143 questions were addressed: (1) how much genetic variation exists among wild populations of *B. oler-*
10
11 144 *acea* in the UK and Spain; (2) how are populations structured in the Atlantic region and how much
12
13 145 differentiation exists between isolated populations; and (3) are there signals of local adaptation to the
14
15 146 environment? The results provide insights into the utility of *B. oleracea* as a crop wild relative genetic
16
17 147 resource for agriculture, as well as shed light on the most likely region of *B. oleracea* domestication.
18
19

20 21 148 **Materials and Methods**

22
23
24 149 Twenty-four populations of *B. oleracea* were chosen from the UK and Spain to cover both a latitudi-
25
26 150 nal and longitudinal gradient of the Atlantic range for genetic analyses (Figure 1i & Table 2). French
27
28 151 populations were not sampled here, but are the focus of a recent genetic analysis by Maggioni *et al.*
29
30 152 (personal communication). Leaves were collected from four individual plants from each population for
31
32 153 DNA extraction, as has been successfully applied to the study of population structure in wild relatives
33
34 154 in the Brassicaceae (Buckley, Holub, Koch, Vergeer, & Mable, 2018). Nazareno, Bemmels, Dick, and
35
36 155 Lohmann (2017) found that compared to “traditional” population genetic markers these smaller sample
37
38 156 sizes are sufficient for various population statistics when large numbers of SNPs are available. The
39
40 157 bedrock for each population was obtained from the British Geological Survey (BGS, 2018) and the Insti-
41
42 158 tuto Geológico y Minero de España (IGME, 2018). The first year a written record of a population exists
43
44 159 was obtained for the UK populations from the Botanical Society of Britain & Ireland (BSBI, 2018). No
45
46 160 equivalent records could be found for the Spanish populations.
47
48

49 161 **Molecular methods**

50
51
52 162 High molecular weight DNA was extracted from the leaves of 96 individuals from 24 populations (Table
53
54 163 2) using DNeasy Plant Mini Kits (QIAGEN, Hilden, Germany) and quantified using a Qubit 2.0 Fluorom-
55
56 164 eter (ThermoFisher Scientific, Waltham, Massachusetts, U.S). Four samples from each population were
57
58 165 sent for library preparation and sequencing at University of Exeter Sequencing Service. Double-digest
59
60

1
2
3
4 166 RADseq libraries were made using a modification of the method in Wu et al. (2016) that allowed Nex-
5
6 167 teraXT indexes (Illumine Corp., USA) to be used for multiplexing samples. In addition, an RYRY spacer
7
8 168 was inserted in the adapter 3' of the Illumina sequencing primer annealing site to provide additional
9
10 169 complexity at the start of read 1 immediately before the Sac1 sticky end. For each sample 400 ng DNA
11
12 170 was fully digested with Sac1 and Mse1 restriction endonucleases and purified using Ampure XP beads.
13
14 171 Illumina compatible i5 adapters were designed to ligate to the at the AGCT-3' sticky end left after Sac1
15
16 172 digest, and Illumina compatible i7 adapters were designed to ligate to the 5'-TA overhangs remaining
17
18 173 after Mse1 digest. Adapter-ligation excess adapters were removed using Ampure XP beads. DNA frag-
19
20 174 ments were amplified by 12 cycles of indexing PCR, purified, size selected (inserts 330-670 bp) and
21
22 175 validated using a Tapestation D1000 HS Screentape (Agilent Technologies Ltd). Libraries were equimo-
23
24 176 lar pooled and the pool concentration was calculated after qPCR. Libraries were denatured, diluted and
25
26 177 sequenced with 125bp paired-end reads on Illumina HiSeq 2500 using SBS High Output reagents v4
27
28 178 (Illumina Corp., USA).

179 **Data processing**

180 Reads were demultiplexed and trimmed to 100 bp using cutadapt (Martin, 2011). These were then
181
182 cleaned and quality filtered using the *process_radtags* pipeline in Stacks v1.47 (Rochette & Catchen,
183
184 2017). Bowtie (v2; Langmead & Salzberg, 2012) and samtools (v1.9; Li et al., 2009) were used to
185
186 align the reads to the *B. oleracea* reference genome (Liu et al., 2014). A catalogue of stacks was then
187
188 created using *ref_map* (Stacks) with the default settings. The *populations* pipeline (Stacks) was used to
189
190 filter the data, and calculate summary statistics. Three datasets were generated with different filtering
191
192 parameters depending on the downstream analysis. Firstly, for dataset 1 (*within individuals*), which was
used to estimate genetic diversity within individuals and in phylogenetic analyses, all individuals were
filtered as a single population, and loci were retained if they had a minimum individual stack depth of
five, a minimum minor allele frequency of 0.01, a maximum observed heterozygosity of 0.7 and were
present in 60% of individuals. Secondly, dataset 2 was generated using the same filtering as dataset 1
but SNPs linked within each RAD locus were avoided by only retaining one SNP at random per locus;
required for population structure analyses (Pritchard, Stephens, & Donnelly, 2000). Finally, for dataset

1
2
3
4
5 193 3 (*within populations*), which was used to calculate genetic distance between populations, individuals
6
7 194 were assigned to their population of origin and loci were retained if present in 50% of the populations.
8
9 195 This filtering was designed to reduce the inclusion of duplicate loci and balance the amount of missing
10
11 196 data with the number of informative loci (Andrews et al., 2016). A minimum stack depth of five is higher
12
13 197 than the default of two, but within the recommended range (Paris, Stevens, & Catchen, 2017), and
14
15 198 helps to remove potential paralogues. Spurious SNPs were avoided by using a minor allele frequency
16
17 199 of > 0.01 (Marandel et al., 2020), and the combination of a maximum observed heterozygosity of 0.7
18
19 200 (70% of the individuals or populations can be heterozygous for each locus) which are present in either
20
21 201 60% of individuals (datasets 1 and 2) or 50% of the populations (dataset 3) retains loci that have been
22
23 202 successfully genotyped across individuals, but are not completely heterozygous. The summary statistics
24
25 203 for each population were calculated in Stacks during the filtering of dataset 3 and included: the number
26
27 204 of private alleles (PRI), expected heterozygosity (H_E), observed heterozygosity (H_O), percentage of
28
29 205 polymorphic loci (%; Table 3), the inbreeding coefficient (F_{IS}) and nucleotide diversity (π ; Supplementary
30
31 206 information).

207 **Data analyses**

32
33
34
35
36 208 Clustering of samples within and between populations was investigated with dataset 1 using RAxML
37
38 209 (v8.2; GTRCAT model and 1000 maximum likelihood bootstrap replicates; Stamatakis, 2014) and visu-
39
40 210 alisation in SplitsTree4 (Huson & Bryant, 2005). To estimate the number of putative genetic clusters
41
42 211 (K) and assess shared genetic ancestry, STRUCTURE (v2.3.4; Pritchard et al., 2000) was used with
43
44 212 dataset 2, so as not to inflate sharing based on multiple SNPs within a RAD locus. A range of K values
45
46 213 were tested (the number of populations successfully sequenced plus one; 1 – 21) using an admixture
47
48 214 model that assumed correlated allele frequencies. For each K, five independent replicates of 100,000
49
50 215 MCMC repetitions, after a burn-in of 10,000 iterations, were run. The most likely K was selected us-
51
52 216 ing the log likelihoods and deltaK (Evanno, Regnaut, & Goudet, 2005). To see if there were significant
53
54 217 differences between estimates of H_E and H_O , pairwise-ANOVAs were carried out in R version 3.4.0 (R
55
56 218 Core Team, 2017) on estimates from dataset 3 based on variant sites alone and all sites. A genetic
57
58 219 distance matrix was created using dataset 3, and the latitude and longitude of each population was
59
60

1
2
3
4 220 used to calculate a geographic distance matrix using 'Haversine' Great Circle Distance in the R package
5
6 221 'geosphere' (Hijmans, 2017). In addition, genetic and geographic matrices were created for Spanish
7
8 222 and UK populations separately, alongside a temporal distance matrix for the year when each population
9
10 223 was first recorded within the UK (first population record; Table 2). Mantel tests were carried out with
11
12 224 9999 replicates on the region-wide matrices and country matrices separately, to assess both the overall
13
14 225 and within country isolation-by-distance. Mantel tests were also carried out on the UK specific matri-
15
16 226 ces to investigate any relationship between the first population records and the genetic and geographic
17
18 227 distances.

19
20 228 A subset of dataset 1 where the soil pH was known was used to investigate the relationship between
21
22 229 soil pH and H_E – e.g., is a higher soil pH associated with lower genetic diversity? A linear model with
23
24 230 soil pH as a predictor variable and H_E as a response variable was run on 21 individuals (across six
25
26 231 populations) from four soil pH classes: Neutral (6.6 - 7.3), Slightly acidic (6.1 - 6.5), Moderately acidic
27
28 232 (5.6 - 6.0) and Strongly acidic (5.0 - 5.5) based on USDA (1998).

29
30 233 In order to identify potential genotype-environment associations, redundancy analyses (RDA) were
31
32 234 carried out using dataset 1 following Forester et al. (2018) with the R packages 'vegan' and 'pysch'
33
34 235 (Oksanen et al., 2017; Reville, 2018). The climate dataset was downloaded from the WorldClim
35
36 236 database at a resolution of 4.5 km (Fick & Hijmans, 2017). This dataset is based on measurements
37
38 237 made between 1970 – 2000. Therefore, it is assumed that any changes in climate will be consistent
39
40 238 enough across the study gradient to maintain differences in the averages and variation between pop-
41
42 239 ulations. The 19 climate variables available from WorldClim for our dataset were checked for pairwise
43
44 240 correlations and the estimated variance inflation factor (VIF). Variables with correlations $> |0.7|$ and
45
46 241 $VIF > 10$ were removed, leaving: 'Annual Mean Temperature', 'Mean Temperature of Wettest Quarter',
47
48 242 'Annual Precipitation' and 'Precipitation Seasonality'. Longitude was included as an additional predictor
49
50 243 variable because it was weakly correlated with climatic variables. Those SNPs that had RDA load-
51
52 244 ings with q-values < 0.1 were considered outlier loci, and were compared to the annotated *B. oleracea*
53
54 245 genome using Bedtools (v.2.17.0; Quinlan & Hall, 2010), followed by a search of the online resource
55
56 246 'Bolbase' (Yu et al., 2013) to investigate putative gene functions.
57
58
59
60

Results

Patterns of genetic diversity

A total of 115,746,909 reads from 76 individuals (20 populations; Table 2) were of sufficient quality and retained for down-stream analysis (average reads per individual: 1,522,986; range: 220,363 – 5,361,799; Supplementary Table 1). For four of the populations, no individuals were successfully sequenced and so these were not included in these analyses. On average 86.3% (range 82.5 - 88.6) of reads mapped to the reference genome (Supplementary Figure 1). Datasets 1 and 2 contained 42,517 and 13,352 SNPs, respectively, across 13,352 RAD-loci (Supplementary Table 2). There were 140,131 SNPs across 53,539 RAD-loci in dataset 3 (Supplementary Information).

Based on variable nucleotide sites only (Table 3), average estimates of genetic diversity (considering H_E) were lower than in the studies cited in Table 1; the average across populations was 0.120 among both UK (range 0.090 – 0.200) and Spanish (range 0.055 – 0.153) populations. Observed heterozygosity was consistently significantly (H_O $p < 0.001$) greater than H_E for all populations and average F_{IS} was similar in the two geographic regions (UK: average = 0.039, range = 0.001 to 0.084; Spain: average = 0.027, range = 0.025 to 0.031). There was thus no evidence of inbreeding (as expected given the genetically controlled self-incompatibility system) but heterozygosity excess was apparent in all populations. The Fortrose population contained 10-fold more private alleles compared to all other populations and had the highest values for both H_E and H_O . Values considering all sites were lower but did not change conclusions about relative patterns of diversity (Table 3).

Population structure

Based on the RAxML tree, the majority of individuals clustered by population, with the exceptions of: (i) two individuals that did not cluster with any population (one in San Juan de Gaxtelugatxe, Spain and one in St Aldehelm's Head, UK), and (ii) an individual from Fortrose (Scotland, UK) that clustered more closely with other Scottish populations than other individuals from Fortrose (Figure 1ii). The most likely number of genetic clusters from STRUCTURE analyses was $K = 12$. Most individuals were admixed, however, six of the UK populations (Fortrose, Auchmithie, Crail, Tynemouth, Whitby and Llantwit Major)

1
2
3
4 273 were dominated by a single genetic ancestry, and two individuals from Fortrose were distinct from both
5
6 274 the third individual from Fortrose and all other samples (Figure 1iii). The dominant genetic ancestry
7
8 275 seen in individuals from Whitby (UK) also dominated the genetic ancestry of individuals from Cabo de
9
10 276 Peñas (Spain), and similarly, the dominant genetic ancestry seen in individuals from Tenby (UK) was
11
12 277 most prevalent in individuals from San Juan de Gaxtelugatxe and Getaría (Spain). There were three
13
14 278 potential regional clusters indicated by the RAxML tree and STRUCTURE analysis: (a) populations
15
16 279 in Scotland; (b) populations closest to the Welsh-English border; and (c) populations in the Basque
17
18 280 Country, Spain (excluding San Sebastian; Figure 1). However, the clustering of populations was not well
19
20 281 resolved and these 'regional clusters' were not always the most geographically close populations (e.g.,
21
22 282 in cluster c, San Sebastian is closer to Getaría than San Juan de Gaxtelugatxe geographically but not
23
24 283 genetically). No isolation-by-distance was predicted by the data either region-wide, or within Spanish
25
26 284 or UK populations alone (Mantel test p-values = 0.474, 0.658 and 0.705, respectively). Furthermore,
27
28 285 no relationship was found between the first record for each of UK populations (Table 2) with either
29
30 286 geographic or genetic distance (Mantel test p-values = 0.114 and 0.933, respectively).
31
32
33

34 287 **Environmental associations**

35
36 288 Overall, environmental variables explained 2.3% (adjusted r-squared) of the variation in the SNPs using
37
38 289 RDA analysis; the strongest association of genotype with the environment was with annual precipitation
39
40 290 (Figure 3). This environmental variation was strong enough to be reflected in the clustering of individ-
41
42 291 uals, including the genetically distinct individuals from Fortrose (UK; Figure 1iii). For example, across
43
44 292 regions, west Scotland and the Basque country experienced the greatest amount of annual precipitation
45
46 293 on average (Figure 2b), whereas the annual mean temperature was greater in the Basque country com-
47
48 294 pared to west Scotland (Figure 2a). Individuals from populations in these regions separated from other
49
50 295 populations in the same direction as annual precipitation, but in opposing directions in relation to annual
51
52 296 mean temperature (Figure 3i). Individuals from Whitby (UK) appear to have experienced a colder, drier
53
54 297 environment than the geographically closest population, Tynemouth (UK), which was also reflected in
55
56 298 the RDA analysis. Linear modelling indicated a non-significant negative trend between genetic diversity
57
58 299 (H_E , H_O & π) and soil pH (i.e. plant genetic diversity decreased as soil pH increased. Only H_E is shown
60

1
2
3
4 but the same relationship was found with H_O & π ; Figure 4).
5

6 There were 2249 unique candidate SNPs associated with the predictor variables from the RDA anal-
7
8 ysis; the majority of these (1039) were most closely associated with 'Mean Temperature of Wettest
9
10 Quarter', followed by 'Precipitation Seasonality' (349), 'Longitude' (333), 'Annual Precipitation' (269) and
11
12 'Annual Mean Temperature' (259). These were fairly evenly distributed across the genome with no indi-
13
14 cations of any single SNP with a large effect. A few SNPs that were more closely associated with annual
15
16 precipitation had strong loadings along axis 1 in the direction of the annual precipitation vector (Figure
17
18 3(ii)). In total, 221 candidate SNPs mapped to unique genes in the *B. oleracea* reference genome, and
19
20 of the top 18, six were annotated as part of the receptor-like kinase family (Table 4).
21
22
23

24 Discussion 25

26
27 The results presented here provide the first genome-wide estimates of genetic variation and population
28
29 genetic structure of wild cabbages collected from across the UK and Spain. Although direct compar-
30
31 isons with cultivated species would be required to rigorously test hypotheses about origins of these
32
33 populations, patterns of variation are consistent with recent linguistic and historical evidence (Maggioni,
34
35 2015; Maggioni et al., 2018) suggesting that the domestication of *B. oleracea* crops occurred in the
36
37 Mediterranean, domesticates were moved by people across Europe, escaped and established wild pop-
38
39 ulations in the Atlantic region. For example, there was no indication of isolation-by-distance from north-
40
41 ern Scotland to Spain ($> 14^\circ$ latitude), which might be expected if these plants were natural colonisers
42
43 following common phylogeographic patterns (e.g., Sharbel, Haubold, & Mitchell-Olds, 2000). Further-
44
45 more, genetic ancestry and clustering analyses suggested that geographically distant populations may
46
47 have similar genetic sources, and could therefore have been established by similar source cultivars.
48
49 The consistent excess of heterozygotes across populations, combined with evidence for admixture from
50
51 STRUCTURE analyses, suggests mixing between 'isolated' populations (Rousset & Raymond, 1995),
52
53 which could be due to interbreeding between cultivated plants growing near the wild populations. This
54
55 highlights the possibility of continued introgression between cultivated and wild plants. Despite the lack
56
57 of geographic genetic population structuring, there were signals of local adaptation to different climates
58
59
60

1
2
3
4 326 based on RDA analyses. In addition, within population genetic diversity estimates were comparable to
5
6 327 other studies (e.g., Christensen et al., 2011; Watson-Jones et al., 2006), and as Watson-Jones et al.
7
8 328 (2006) found, lower genetic diversity estimates were associated with higher soil pH. Therefore, these
9
10 329 wild populations could hold useful adaptive alleles for plant breeding, and a suitable approach to investi-
11
12 330 gate traits of agricultural interest (e.g., drought tolerance) could be to choose populations based on their
13
14 331 environment of origin. However, further sequencing of a range of cultivars from different geographic
15
16 332 regions would be required to further test these hypotheses.

333 **Patterns of Genetic Diversity**

334 Although the magnitude of estimates of genetic diversity based on the ddRADseq data presented here
335 were lower than in previous studies (see Table 1) using allozymes (Lanner-Herrera, Gustafeson, Filt, &
336 Bryngelsson, 1996; Lázaro & Aguinagalde, 1998; Raybould et al., 1999), microsatellites (Raybould et
337 al., 1999) or AFLPs (Watson-Jones et al., 2006; Christensen et al., 2011), patterns of variation within
338 the UK and Spain were strikingly similar to one another. Most populations also showed a relatively con-
339 sistent excess of heterozygosity. These similarities could provide evidence for relatively recent origins
340 of populations in the two regions, but whether this was from feralisation of cultivars or natural differ-
341 entiation after natural colonisation cannot be distinguished by the data. Although there has been an
342 ongoing debate as to the origin of wild *B. oleracea* populations in the Atlantic region (Song, Osborn,
343 & Williams, 1990; Allender et al., 2007; Maggioni, 2015), domestication of *B. oleracea* in the Mediter-
344 ranean region has been suggested by other genetic, phenotypic and linguistic studies (Mitchell, 1976;
345 Maggioni, 2015; Maggioni et al., 2018). The subsequent movement of *B. oleracea* cultivars across Eu-
346 rope could then have resulted in a much narrower bottleneck than the initial domestication bottleneck
347 in the Mediterranean as it removed the chance of gene flow from the wild relatives they originated from
348 (Kofsky, Zhang, & Song, 2018). Consistent with this hypothesis, although the putative Mediterranean
349 progenitor species remains unknown, Allender et al. (2007) found much greater estimates of genetic
350 diversity within potential progenitor species from the Mediterranean region than either previous genetic
351 diversity estimates made in *B. oleracea* (e.g., Christensen et al., 2011; Watson-Jones et al., 2006) or in
352 this study.

Population structure

Several of the analyses here suggest less population structuring than might be expected in such geographically distinct populations if natural range expansion followed by isolation occurred. In this dataset, since the first recorded population (Tenby in 1773), one to three new populations have been recorded every thirty years within the UK (Table 2). However, neither the date the UK populations were first recorded, nor the genetic distances between populations in the UK and Spain, had a geographical pattern (i.e. no isolation by distance). Furthermore, although the majority of individuals clustered by population and some regional clustering was seen (Figure 1), it would not be possible to predict whether two individuals from geographically close or geographically distant populations are more genetically similar to each other. For example, Fowey and Prussia Cove (UK populations), and West Looe and Cabo de Peñas (UK and Spanish populations respectively), clustered together and shared more genetic ancestry than Fowey and West Looe, which are the closest geographically. Although more sampling would be required to explicitly test it, the evidence here suggests that these plants have not colonised the Atlantic region following common phylogeographic patterns (e.g., Sharbel et al., 2000) and therefore is consistent with *B. oleracea* domestication occurring outside of the Atlantic region. This is in line with results from other genetic, phenotypic and linguistic studies, which suggest the Mediterranean region is the most likely location for *B. oleracea* domestication (Maggioni, 2015; Maggioni et al., 2018; Mitchell, 1976).

The genetic ancestry and clustering analyses hint that populations could have been established by escapees from different cultivars. The majority of individuals were assigned to multiple sources of genetic ancestry (Figure 1iii), however, there were also cases where one putative source dominated at the individual- and population-levels, which could be the overall genetic background from the original source cultivar. Interestingly, there were two distinct individuals from Fortrose (10-fold more private alleles than other populations; Table 3) with a source that was assigned to no other individuals. Due to the ease of interbreeding between cultivars (Allender et al., 2007), this could indicate that these two Fortrose individuals are recent escapees from a different source population (e.g., local gardens), which are yet to have mixed with other individuals within the population. Furthermore, the more recent record of the population at Fortrose (1968), and the lack of assignment to other populations, suggests that this genetic background could be from a cultivar that has not been grown for a long period of time or widely

1
2
3
4 381 around the Atlantic coastlines. The excess of heterozygotes (H_O was significantly greater than H_E) and
5
6 382 the general mix of shared genetic ancestry across such a wide geographical area in distinct populations,
7
8 383 could also be an indication of continued introgression into these wild populations from agricultural and
9
10 384 horticultural sources. It would be interesting to identify popular cultivars in the local areas of these
11
12 385 populations, including any changes in the preferred cultivars through time, to investigate patterns of
13
14 386 introgression in more detail. Such direct comparisons with cultivars could identify the most likely founder
15
16 387 of these populations.

17
18 388 Using chloroplast microsatellite DNA markers, Allender et al. (2007) found two haplotypes in *B.*
19
20 389 *oleracea* around the coasts of the UK; out of sixteen populations, fourteen were C:01 and two were
21
22 390 C:04. The two populations with the C:04 haplotype were in Tyne & Wear, in the northeast of England;
23
24 391 in the current study, this area is represented by the Tynemouth and Whitby populations. In line with the
25
26 392 rarity of the chloroplast haplotypes identified in this region in the previous study, these two populations
27
28 393 clustered most closely with populations not sampled by Allender et al. (2007); Tynemouth clustered with
29
30 394 Fortrose, Scotland, and Whitby with the Spanish population Cabo de Peñas. Based on this information,
31
32 395 it might be expected that the chloroplast haplotypes of Fortrose and Cabo de Peñas would also be
33
34 396 C:04. In addition, the C:01 haplotype found in the majority of the UK populations was also found in
35
36 397 four other species of *Brassica* (Allender et al., 2007), suggesting either that this is the ancestral form or
37
38 398 introgression between species. A combination of nuclear and chloroplast information could be useful for
39
40 399 disentangling the population histories further, particularly in relation to identifying introgression.

41
42 400 Knowledge of the founding cultivars would be useful for both plant breeders and those interested
43
44 401 in invasive species. It could provide insights into how different cultivars have adapted (and therefore
45
46 402 may adapt in the future) to different environmental conditions, and could also be thought of as a way
47
48 403 to compare invasion success within a species. *Brassica oleracea* lack the characteristics thought to be
49
50 404 fundamental for establishment in novel locations (invasions; Funk et al., 2016), but perhaps amongst
51
52 405 the huge phenotypic variation found within this species, some traits are more likely to lead to successful
53
54 406 'invasions' of particular cultivars compared to others. For example, a cultivated Danish kale was the
55
56 407 most likely source for a wild population found in Denmark (based on AFLP markers; Christensen et
57
58 408 al., 2011), and it could be that all the Atlantic populations were established by different kale cultivars.

1
2
3
4
5 409 Overall, populations of *B. oleracea* growing along Atlantic coasts would be an excellent study system to
6
7 410 improve understanding of invasive species that are likely to harbour useful adaptive traits for agriculture.

8
9 411 While comparisons with published whole genome sequence data or other types of genotype by
10
11 412 sequencing approaches (e.g., Stansell et al., 2018) for cultivated *B. oleracea* would be interesting to
12
13 413 more explicitly test origins of the populations studied here, there are several issues with ddRAD data
14
15 414 that would make this challenging and potentially hard to interpret. A benefit of ddRAD sequencing is
16
17 415 the generation of discrete loci that are standardised to the same length. However, the resulting short
18
19 416 sequence segments normally contain only one or a few SNPs, which does not allow accurate assign-
20
21 417 ment of paralogs in highly duplicated and rearranged genomes such as found in the Brassicaceae (e.g.,
22
23 418 Schranz, Lysak, & Mitchell-Olds, 2006). Instead, filtering pipelines to allow population genetics analyses
24
25 419 based on ddRAD data are designed to be conservative (Paris et al., 2017; Marandel et al., 2020). This
26
27 420 filtering results in fewer loci retained, but it should reduce risks of including duplicates. In the current
28
29 421 study, excess heterozygosity was observed consistently across populations, which could suggest his-
30
31 422 torical introgression. Although we cannot completely rule out the influence of combining duplicates (Ilut,
32
33 423 Nydam, & Hare, 2014), the highly consistent patterns of excess suggest that all populations would have
34
35 424 been affected similarly, enabling interpretations of relative variation within and between populations. The
36
37 425 admixture suggested by the STRUCTURE analyses also supports the role of introgression in the histo-
38
39 426 ries of the studied populations. However, mapping of the ddRAD reads to multiple reference genomes
40
41 427 or to data generated based on different restriction enzymes would be more problematic.

42 43 44 428 **Environmental associations**

45
46 429 Despite the general lack of geographic clustering, there was evidence of local adaptation to the vary-
47
48 430 ing environments using redundancy analyses, particularly to annual precipitation (Figure 3). Although
49
50 431 Watson-Jones et al. (2006) found some population structuring within the UK, the same result was not
51
52 432 found in this study (i.e. no isolation-by-distance within the UK). Furthermore, no evidence of population
53
54 433 structuring was found in the Spanish populations here, and Maggioni *et al.* (personal communication)
55
56 434 found no evidence of population structuring in French Atlantic populations. These results could also be
57
58 435 correlated with annual precipitation; perhaps the strong variation in annual precipitation in the UK (e.g.,
59
60

1
2
3
4 436 a strong west-east gradient) is causing more differentiation between these populations, whereas along
5
6 437 the French range annual precipitation has a smaller gradient. One reason for the importance of annual
7
8 438 precipitation other than water availability could be the influence of precipitation on soil pH. Soil pH is
9
10 439 primarily determined by bedrock, but is also altered by precipitation through leaching of compounds
11
12 440 such as calcium carbonate (Kinzel, 1983). Therefore, although slightly alkaline to neutral soils tend to
13
14 441 form over limestone, secondary acidification can occur under higher precipitation regimes. The soil pH
15
16 442 values recorded here ranged from neutral to strongly acidic (Figure 4). Furthermore, the bedrock of a
17
18 443 large proportion of the populations used here (Table 2) differ from the limestone and chalk cliffs that wild
19
20 444 *B. oleracea* are thought to be predominantly found on (Christensen et al., 2011). For those individuals
21
22 445 where the soil pH was known, the same trend was found here as by Watson-Jones et al. (2006), with
23
24 446 a decrease in plant genetic diversity as soil pH increased (Figure 4). For agriculture and horticulture,
25
26 447 soil pH is an important consideration (Tilman, Balzer, Hill, & Befort, 2011). The change in plant genetic
27
28 448 diversity suggests that soil pH is a strong selective pressure in the wild, causing an adaptive ecolog-
29
30 449 ical bottleneck in locations where it is higher, resulting in lower genetic diversity. These indications of
31
32 450 local adaptation despite a lack of population structure highlight environmental variables that could be
33
34 451 investigated further in wild populations of *B. oleracea*, which regardless of their origin are surviving.

35
36 452 Alongside survival, a huge concern for food security related to climate change is the ability of crop
37
38 453 plants to remain productive under rapidly changing environmental conditions (Lasky et al., 2015). Ob-
39
40 454 taining accurate phenotypic data for adaptive traits is a major barrier as we often do not know the com-
41
42 455 bination of traits that underlie differences in fitness or how these vary with the environment (Kooyers,
43
44 456 Greenlee, Colicchio, Oh, & Blackman, 2015). Although some traits will be locally adaptive due to large
45
46 457 effect loci, the vast majority of adaptive traits are likely to have a polygenic basis (Rockman, 2012),
47
48 458 particularly in the case of multi-trait phenotypes related to environmental gradients. Our results match
49
50 459 these expectations, as no large effect loci were found; however, some were more significantly associ-
51
52 460 ated with the assessed environmental variation than others. The most likely assignment for six of the
53
54 461 top 18 candidate genes was to the receptor-like kinase family (Table 4). This gene family underwent an
55
56 462 expansion that is believed to be a plant-specific adaptation for pathogen defence (Afzal, Wood, & Light-
57
58 463 foot, 2008). Interestingly, Zhang et al. (2014) also found differences in genes related to plant defence
59
60

1
2
3
4 464 when investigating adaptations of rice (*Oryza* sp.) across four continents. These results highlight the
5
6 465 fundamental importance of the immune system to fitness, and suggest that it could be related to envi-
7
8 466 ronmental differences across different spatial scales. Given that immune system genes are among the
9
10 467 best candidates for local adaptation, there is a potential connection between plant genetic diversity, soil
11
12 468 pH and pathogens. It would be interesting to investigate whether less acidic soils host more pathogens,
13
14 469 increasing the selective pressure on the plants and decreasing the plant genetic diversity in these soils.
15
16 470 Overall, the impact of climate change on the spread of virulence of plant pathogens and herbivores, and
17
18 471 the phenological mismatches that may occur between interacting species remain unknown (De Lucia,
19
20 472 Nability, Zavala, & Berenbaum, 2012; Fisher et al., 2012; Yang & Rudolf, 2010). What is clear is that plant
21
22 473 defence will continue to be an important component of crop productivity, warranting further research.

23
24 474 Overall, the results presented here supported the hypothesis that wild populations of *B. oleracea* in
25
26 475 the Atlantic region were established by plants from agricultural and/or horticultural sources. In addition,
27
28 476 regardless of their origin, these wild populations are likely to contain useful genetic resources and should
29
30 477 be considered as valuable populations of a crop wild relative to be investigated further.

31 32 33 34 478 **Acknowledgements**

35
36
37 479 EAM was funded by a University of Glasgow Lord Kelvin Adam Smith PhD studentship; UZI was
38
39 480 funded by a NERC Independent Research Fellowship (NE/L011956); CAC is supported by the BBSRC
40
41 481 (BB/P004202/1); KAM utilised equipment funded by the Wellcome Trust Institutional Strategic Support
42
43 482 Fund (WT097835MF), Wellcome Trust Multi User Equipment Award (WT101650MA) and BBSRC LOLA
44
45 483 award (BB/K003240/1). Part of the work was supported by a British Society for Plant Pathology summer
46
47 484 studentship, and grants from the Botanical Research Fund, and the Blodwen Lloyd Bins trust funded
48
49 485 through the Glasgow Natural History Society. None of the sponsors had any role in the study design,
50
51 486 data collection, analysis, and interpretation or any aspects during the write up and publication of this
52
53 487 work. We thank anonymous reviewers for helpful comments, Danijela Dimitrijević and Deborah Davy for
54
55 488 assistance in the field, and Dr Lorenzo Maggioni for useful discussions.
56
57
58
59
60

Conflict of Interest

None declared.

References

- Afzal, A. J., Wood, A. J., & Lightfoot, D. A. (2008). Plant receptor-like serine threonine kinases: roles in signaling and plant defense. *Molecular Plant-Microbe Interactions*, *21*(5), 507–517.
- Allender, C., Allainguillaume, J., Lynn, J., & King, G. J. (2007). Simple sequence repeats reveal uneven distribution of genetic diversity in chloroplast genomes of *Brassica oleracea* L. and (n= 9) wild relatives. *Theoretical and Applied Genetics*, *114*(4), 609–618.
- Andrews, K. R., Good, J. M., Miller, M. R., Luikart, G., & Hohenlohe, P. A. (2016). Harnessing the power of RADseq for ecological and evolutionary genomics. *Nature Reviews Genetics*, *17*(2), 81–92.
- BGS. (2018). *British Geological Survey*. <http://www.bgs.ac.uk/>. (Accessed: 2018-08-26)
- BSBI. (2018). *Botanical Society of Britain & Ireland*. <https://bsbi.org>. (Accessed: 2018-08-26)
- Buckley, J., Holub, E. B., Koch, M. A., Vergeer, P., & Mable, B. K. (2018). Restriction associated DNA-genotyping at multiple spatial scales in *Arabidopsis lyrata* reveals signatures of pathogen-mediated selection. *BMC Genomics*, *19*(1), 496.
- Charlesworth, D., Vekemans, X., Castric, V., & Glémin, S. (2005). Plant self-incompatibility systems: a molecular evolutionary perspective. *New Phytologist*, *168*(1), 61–69.
- Christensen, S., von Bothmer, R., Poulsen, G., Maggioni, L., Phillip, M., Andersen, B. A., & Jørgensen, R. B. (2011). AFLP analysis of genetic diversity in leafy kale (*Brassica oleracea* L. convar. *acephala* (DC.) Alef.) landraces, cultivars and wild populations in Europe. *Genetic Resources and Crop Evolution*, *58*(5), 657–666.
- Conner, J. K. (2002). Genetic mechanisms of floral trait correlations in a natural population. *Nature*, *420*(6914), 407–410.
- Darwin, C. R. (1859). *The origin of species*. London: John Murray.
- De Lucia, E., Nabity, P., Zavala, J., & Berenbaum, M. (2012). Climate change: resetting plant-insect interactions. *Plant Physiology*, pp–112.
- Diamond, J. (2002). Evolution, consequences and future of plant and animal domestication. *Nature*, *418*(6898), 700.
- Euro+Med PlantBase. (2020). <http://ww2.bgbm.org/EuroPlusMed/>. (Accessed: 2020-05-23)
- Evanno, G., Regnaut, S., & Goudet, J. (2005). Detecting the number of clusters of individuals using the software STRUCTURE: a simulation study. *Molecular Ecology*, *14*(8), 2611–2620.
- Falconer, D. S., & Mackay, T. F. C. (1996). *Introduction to Quantitative Genetics* (4th ed.). Harlow, UK: Longman.
- Fick, S. E., & Hijmans, R. J. (2017). WorldClim 2: new 1-km spatial resolution climate surfaces for global land areas. *International Journal of Climatology*, *37*(12), 4302–4315.
- Fisher, M. C., Henk, D. A., Briggs, C. J., Brownstein, J. S., Madoff, L. C., McCraw, S. L., & Gurr, S. J. (2012). Emerging fungal threats to animal, plant and ecosystem health. *Nature*, *484*(7393), 186.
- Forester, B. R., Lasky, J. R., Wagner, H. H., & Urban, D. L. (2018). Comparing methods for detecting multilocus adaptation with multivariate genotype–environment associations. *Molecular Ecology*, *27*(9), 2215–2233.
- Funk, J. L., Standish, R. J., Stock, W. D., & Valladares, F. (2016). Plant functional traits of dominant native and invasive species in mediterranean-climate ecosystems. *Ecology*, *97*(1), 75–83.
- Gao, Y., Yin, S., Wu, L., Dai, D., Wang, H., Liu, C., & Tang, L. (2017). Genetic diversity and structure of wild and cultivated *Amorphophallus paeoniifolius* populations in southwestern China as revealed by RAD-seq. *Scientific Reports*, *7*(1), 14183.
- Gómez-Campo, C., & Prakash, S. (1999). Origin and domestication. In *Developments in plant genetics and breeding* (Vol. 4, pp. 33–58). Elsevier.
- Henriksen, R., Gering, E., & Wright, D. (2018). Feralisation—The understudied counterpoint to domestication. In *Origin and evolution of biodiversity* (pp. 183–195). Springer.
- Hijmans, R. J. (2017). geosphere: Spherical Trigonometry [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=geosphere> (R package version 1.5-7)

- 1
2
3
4
540 Hoisington, D., Khairallah, M., Reeves, T., Ribaut, J.-M., Skovmand, B., Taba, S., & Warburton, M.
541 (1999). Plant genetic resources: What can they contribute toward increased crop productivity?
542 *Proceedings of the National Academy of Sciences*, 96(11), 5937–5943.
- 543 Huson, D. H., & Bryant, D. (2005). Application of phylogenetic networks in evolutionary studies. *Molec-*
544 *ular Biology and Evolution*, 23(2), 254–267.
- 545 IGME. (2018). *Instituto Geológico y Minero de España*. <http://www.igme.es>. (Accessed: 2018-08-26)
- 546 Ilut, D. C., Nydam, M. L., & Hare, M. P. (2014). Defining loci in restriction-based reduced representation
547 genomic data from nonmodel species: sources of bias and diagnostics for optimal clustering.
548 *BioMed Research International*, 2014.
- 549 Kinzel, H. (1983). Influence of limestone, silicates and soil pH on vegetation. In *Physiological plant*
550 *ecology iii* (pp. 201–244). Springer.
- 551 Kitashiba, H., & Nasrallah, J. B. (2014). Self-incompatibility in Brassicaceae crops: lessons for interspe-
552 cific incompatibility. *Breeding Science*, 64(1), 23–37.
- 553 Kofsky, J., Zhang, H., & Song, B.-H. (2018). The untapped genetic reservoir: The past, current, and
554 future applications of the wild soybean (*Glycine soja*). *Frontiers in Plant Science*, 9.
- 555 Kooyers, N. J., Greenlee, A. B., Colicchio, J. M., Oh, M., & Blackman, B. K. (2015). Replicate altitudinal
556 clines reveal that evolutionary flexibility underlies adaptation to drought stress in annual *Mimulus*
557 *guttatus*. *New Phytologist*, 206(1), 152–165.
- 558 Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods*,
559 9(4), 357.
- 560 Lanner-Herrera, C., Gustafeson, M., Filt, A., & Bryngelsson, T. (1996). Diversity in natural populations
561 of wild *Brassica oleracea* as estimated by isozyme and RAPD analysis. *Genetic Resources and*
562 *Crop Evolution*, 43(1), 13–23.
- 563 Lasky, J. R., Upadhyaya, H. D., Ramu, P., Deshpande, S., Hash, C. T., Bonnette, J., ... others (2015).
564 Genome-environment associations in sorghum landraces predict adaptive traits. *Science Ad-*
565 *vances*, 1(6), e1400218.
- 566 Lázaro, A., & Aguinalgalde, I. (1998). Genetic diversity in *Brassica oleracea* L.(Cruciferae) and wild
567 relatives (2 n= 18) using isozymes. *Annals of Botany*, 82(6), 821–828.
- 568 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., ... Durbin, R. (2009). The sequence
569 alignment/map format and SAMtools. *Bioinformatics*, 25(16), 2078–2079.
- 570 Liu, S., Liu, Y., Yang, X., Tong, C., Edwards, D., Parkin, I. A., ... others (2014). The *Brassica oleracea*
571 genome reveals the asymmetrical evolution of polyploid genomes. *Nature Communications*, 5.
- 572 Maggioni, L. (2015). *Domestication of Brassica oleracea l.* (Unpublished doctoral dissertation). Swedish
573 University of Agricultural Studies, Alnarp, Sweden.
- 574 Maggioni, L., von Bothmer, R., Poulsen, G., & Lipman, E. (2018). Domestication, diversity and use
575 of *Brassica oleracea* L., based on ancient Greek and Latin texts. *Genetic Resources and Crop*
576 *Evolution*, 65(1), 137–159.
- 577 Marandel, F., Charrier, G., Lamy, J.-B., Le Cam, S., Lorange, P., & Trenkel, V. M. (2020). Estimating
578 effective population size using RADseq: Effects of SNP selection and sample size. *Ecology and*
579 *Evolution*, 10(4), 1929–1937.
- 580 Martin, M. (2011, May). Cutadapt removes adapter sequences from high-throughput sequencing
581 reads. *EMBnet journal*, 17(1), 10–12. Retrieved from [http://journal.embnet.org/index.php/](http://journal.embnet.org/index.php/embnetjournal/article/view/200)
582 [embnetjournal/article/view/200](http://journal.embnet.org/index.php/embnetjournal/article/view/200) doi: <http://dx.doi.org/10.14806/ej.17.1.200>
- 583 Milla, R., Osborne, C. P., Turcotte, M. M., & Violle, C. (2015). Plant domestication through an ecological
584 lens. *Trends in Ecology & Evolution*.
- 585 Mitchell, N. (1976). The Status of *Brassica oleracea* L. subsp. *oleracea* (wild cabbage) in the British
586 Isles. *Watsonia*, 11, 97–103.
- 587 Nazareno, A. G., Bemmels, J. B., Dick, C. W., & Lohmann, L. G. (2017). Minimum sample sizes for
588 population genomics: An empirical study from an Amazonian plant species. *Molecular Ecology*
589 *Resources*, 17(6), 1136–1147.
- 590 Nei, M. (1973). Analysis of gene diversity in subdivided populations. *Proceedings of the National*
591 *Academy of Sciences*, 70(12), 3321–3323.
- 592 Oksanen, J., Blanchet, F. G., Friendly, M., Kindt, R., Legendre, P., McGlinn, D., ... Wagner, H. (2017).
593 vegan: Community Ecology Package [Computer software manual]. Retrieved from [https://CRAN](https://CRAN.R-project.org/package=vegan)
594 [.R-project.org/package=vegan](https://CRAN.R-project.org/package=vegan) (R package version 2.4-3)

- 1
2
3
4
595 Paris, J. R., Stevens, J. R., & Catchen, J. M. (2017). Lost in parameter space: a road map for stacks.
596 *Methods in Ecology and Evolution*, 8(10), 1360–1373.
- 597 Prentis, P. J., Wilson, J. R., Dormontt, E. E., Richardson, D. M., & Lowe, A. J. (2008). Adaptive evolution
598 in invasive species. *Trends in Plant Science*, 13(6), 288–294.
- 599 Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of population structure using multilocus
600 genotype data. *Genetics*, 155(2), 945–959.
- 601 Purugganan, M. D., & Fuller, D. Q. (2009). The nature of selection during plant domestication. *Nature*,
602 457(7231), 843–848.
- 603 Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features.
604 *Bioinformatics*, 26(6), 841–842.
- 605 R Core Team. (2017). R: A language and environment for statistical computing [Computer software
606 manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- 607 Rauf, S., Teixeira da Silva, J., Khan, A. A., & Naveed, A. (2010). Consequences of plant breeding on
608 genetic diversity. *International Journal of Plant Breeding*, 4(1), 1–21.
- 609 Raybould, A., Mogg, R., Clarke, R., Gliddon, C., & Gray, A. (1999). Variation and population structure
610 at microsatellite and isozyme loci in wild cabbage (*Brassica oleracea* L.) in Dorset (UK). *Genetic
611 Resources and Crop Evolution*, 46(4), 351–360.
- 612 Revelle, W. (2018). psych: Procedures for Psychological, Psychometric, and Personality Research
613 [Computer software manual]. Evanston, Illinois. Retrieved from [https://CRAN.R-project.org/
614 package=psych](https://CRAN.R-project.org/package=psych) (R package version 1.8.10)
- 615 Rochette, N. C., & Catchen, J. M. (2017). Deriving genotypes from RAD-seq short-read data using
616 Stacks. *Nature Protocols*, 12(12), 2640.
- 617 Rockman, M. V. (2012). The QTN program and the alleles that matter for evolution: all that's gold does
618 not glitter. *Evolution*, 66(1), 1–17.
- 619 Rousset, F., & Raymond, M. (1995). Testing heterozygote excess and deficiency. *Genetics*, 140(4),
620 1413–1419.
- 621 Sánchez-Yélamo, M. (2014). Characterisation of wild cabbage (*Brassica oleracea* L.) based on isoen-
622 zyme data. considerations on the current status of this taxon in Spain. *Genetic Resources and
623 Crop Evolution*, 61(7), 1295–1306.
- 624 Scheepens, J., Frei, E. S., Armbruster, G. F., & Stöcklin, J. (2012). Pollen dispersal and gene flow within
625 and into a population of the alpine monocarpic plant *Campanula thyrsoidea*. *Annals of Botany*,
626 110(7), 1479–1488.
- 627 Schranz, M. E., Lysak, M. A., & Mitchell-Olds, T. (2006). The ABC's of comparative genomics in the
628 Brassicaceae: building blocks of crucifer genomes. *Trends in Plant Science*, 11(11), 535–542.
- 629 Sharbel, T. F., Haubold, B., & Mitchell-Olds, T. (2000). Genetic isolation by distance in *Arabidopsis
630 thaliana*: biogeography and postglacial colonization of Europe. *Molecular Ecology*, 9(12), 2109–
631 2118.
- 632 Slatkin, M. (1987). Gene flow and the geographic structure of natural populations. *Science*, 236(4803),
633 787–792.
- 634 Song, K., Osborn, T. C., & Williams, P. H. (1990). *Brassica* taxonomy based on nuclear restriction
635 fragment length polymorphisms (RFLPs). *Theoretical and Applied Genetics*, 79(4), 497–506.
- 636 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
637 phylogenies. *Bioinformatics*, 30(9), 1312–1313.
- 638 Stansell, Z., Hyma, K., Fresnedo-Ramírez, J., Sun, Q., Mitchell, S., Björkman, T., & Hua, J. (2018).
639 Genotyping-by-sequencing of *Brassica oleracea* vegetables reveals unique phylogenetic patterns,
640 population structure and domestication footprints. *Horticulture Research*, 5.
- 641 Tilman, D., Balzer, C., Hill, J., & Befort, B. L. (2011). Global food demand and the sustainable intensifi-
642 cation of agriculture. *Proceedings of the National Academy of Sciences*, 108(50), 20260–20264.
- 643 USDA, D. o. A. (1998). *Soil Quality Indicators: pH*. [https://www.nrcs.usda.gov/Internet/
644 FSE_DOCUMENTS/nrcs142p2_052208.pdf](https://www.nrcs.usda.gov/Internet/FSE_DOCUMENTS/nrcs142p2_052208.pdf). (Accessed: 2016-04-20)
- 645 von Wettberg, E. J., Chang, P. L., Başdemir, F., Carrasquilla-Garcia, N., Korbu, L. B., Moenga, S. M., . . .
646 others (2018). Ecology and genomics of an important crop wild relative as a prelude to agricultural
647 innovation. *Nature Communications*, 9(1), 649.
- 648 Walley, P. G., Teakle, G. R., Moore, J. D., Allender, C. J., Pink, D. A., Buchanan-Wollaston, V., & Barker,
649 G. C. (2012). Developing genetic resources for pre-breeding in *Brassica oleracea* L.: an overview
650 of the UK perspective. *Journal of Plant Biotechnology*, 39(1), 62–68.

- 651 Watson-Jones, S., Maxted, N., & Ford-Lloyd, B. (2006). Population baseline data for monitoring genetic
652 diversity loss for 2010: a case study for *Brassica* species in the UK. *Biological Conservation*,
653 132(4), 490–499.
- 654 Wichmann, M. C., Alexander, M. J., Soons, M. B., Galsworthy, S., Dunne, L., Gould, R., ... Bullock,
655 J. M. (2009). Human-mediated dispersal of seeds over long distances. *Proceedings of the Royal
656 Society of London B: Biological Sciences*, 276(1656), 523–532.
- 657 Wu, Z., Wang, B., Chen, X., Wu, J., King, G. J., Xiao, Y., & Liu, K. (2016). Evaluation of linkage
658 disequilibrium pattern and association study on seed oil content in *Brassica napus* using ddRAD
659 sequencing. *PLoS One*, 11(1), e0146383.
- 660 Xiao, Z., Han, F., Hu, Y., Xue, Y., Fang, Z., Yang, L., ... others (2019). Overcoming cabbage crossing
661 incompatibility by the development and application of self-compatibility-QTL-specific markers and
662 genome-wide background analysis. *Frontiers in Plant Science*, 10, 189.
- 663 Yang, L. H., & Rudolf, V. (2010). Phenology, ontogeny and the effects of climate change on the timing
664 of species interactions. *Ecology Letters*, 13(1), 1–10.
- 665 Yousef, E. A., Mueller, T., Börner, A., & Schmid, K. J. (2018). Comparative analysis of genetic diver-
666 sity and differentiation of cauliflower (*Brassica oleracea* var. botrytis) accessions from two *ex situ*
667 genebanks. *PLoS One*, 13(2), e0192062.
- 668 Yu, J., Zhao, M., Wang, X., Tong, C., Huang, S., Tehrim, S., ... Liu, S. (2013). Bolbase: a comprehensive
669 genomics database for *Brassica oleracea*. *BMC Genomics*, 14(1), 664.
- 670 Zhang, Q.-J., Zhu, T., Xia, E.-H., Shi, C., Liu, Y.-L., Zhang, Y., ... others (2014). Rapid diversification of
671 five *Oryza* AA genomes associated with rice adaptation. *Proceedings of the National Academy of
672 Sciences*, 111(46), E4954–E4962.

673 **Data Accessibility**

674 The sequencing data and associated meta data are available on the European Nucleotide Archive under
675 the study accession number: PRJEB38464 (<http://www.ebi.ac.uk/ena/data/view/PRJEB38464>)

676 **Author Contributions**

677 EAM and BKM conceived and proposed the study to the co-authors; EAM, BKM, CAC, and UZI then
678 developed the study design. BKM, CAC, and UZI obtained the main funding award. EAM and EAK
679 collected the data. EAM carried out statistical analysis, and wrote the paper. KAM carried out the
680 double-digest RAD-sequencing. All authors reviewed the final version of the manuscript and agreed to
681 its content before submission.

Table 1: Estimates of genetic diversity within wild *B. oleracea* populations from previous studies using different molecular markers. H_E is expected heterozygosity estimated using Nei's gene diversity (Nei, 1973).

Study	Molecular marker	H_E	Populations
Lanner-Herrera et al. (1996)	Isozymes	0.10 – 0.56	France, Spain, UK
Lázaro and Aguinagalde (1998)	Isozymes	0.26 – 0.30	France, Spain, UK
Raybould et al. (1999)	Isozymes	0.40 (0.18 – 0.41)†	UK
Raybould et al. (1999)	Microsatellites	0.36 (0.21 – 0.33)†	UK
Watson-Jones et al. (2006)	AFLPs	0.19 – 0.33	UK
Christensen et al. (2011)	AFLPs	0.23, 0.20	Spain, UK
Maggioni <i>et al.</i> (pers. comm. 2019)	AFLPs	0.25	France

† – pooled population H_E with the range of estimates from individual populations shown in brackets.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

24

Table 2: A summary of the natural populations of *B. oleracea* used in this study, including: the bedrock, the first time the population was recorded, the number of individuals sequenced, and the number of individuals included in down-stream analyses.

Region	Population	Bedrock†	First population record‡	Number sequenced	Number included§
ES	Auchmithie	Red basic sandstone	1913	4	4
ES	Crail	Sandstone & limestone	1840	4	4
ES	Fortrose	Sandstone	1968	4	3
WS	Kildonan Castle	Sandstone & limestone	1987	4	4
NEE	Tynemouth	Sandstone & limestone	1805	4	4
NEE	Staithe	Shale & sandstone	1831	4	–
NEE	Whitby	Shale	1906	4	4
NW	Little Orme	Limestone	1895	4	–
SW	Tenby	Siltstone & sandstone	1773	4	4
SW	Llantwit Major	Limestone	1850	4	4
SWE	Prussia Cove	Slate, shale & siltstone	1871	4	4
SWE	Fowey	Shale & siltstone	1805	4	4
SWE	West Looe	Siltstone & sandstone	1971	4	2
SWE	St. Aldhelm's Head	Limestone	1933	4	4
A	Cudillero	Slate & sandstone		4	4
A	Playa de Xágo	Sandstone & Dolomite		4	4
A	Cabo de Peñas	Slate & quartzite		4	4
A	Playas de Viodo	Slate & shale		4	4
A	Tazonas	Dolomite & limestone		4	–
C	Playa Pedrero	Quartzites		4	4
C	La Franca	Quartzites		4	–
BC	San Juan de Gaxtelugatxe	Limestone		4	4
BC	Getaría	Limestone		4	4
BC	San Sebastian	Calcareous sandstone		4	3
Total:				96	76

† data obtained from the British Geological Survey (<https://www.bgs.ac.uk>) and the Instituto Geológico y Minero de España (<http://www.igme.es>). Region codes: ES – East Scotland, WS – West Scotland, NEE – North-eastern England, NW – North Wales, SW – South Wales, SWE – South-western England, A – Asturias Spain, C – Cantabrica Spain, BC – Basque Country Spain. ‡ data obtained from the Botanical Society of Britain & Ireland (<https://bsbi.org>). § indicates where data was lost in quality filtering of sequences and not included in down-stream analyses.

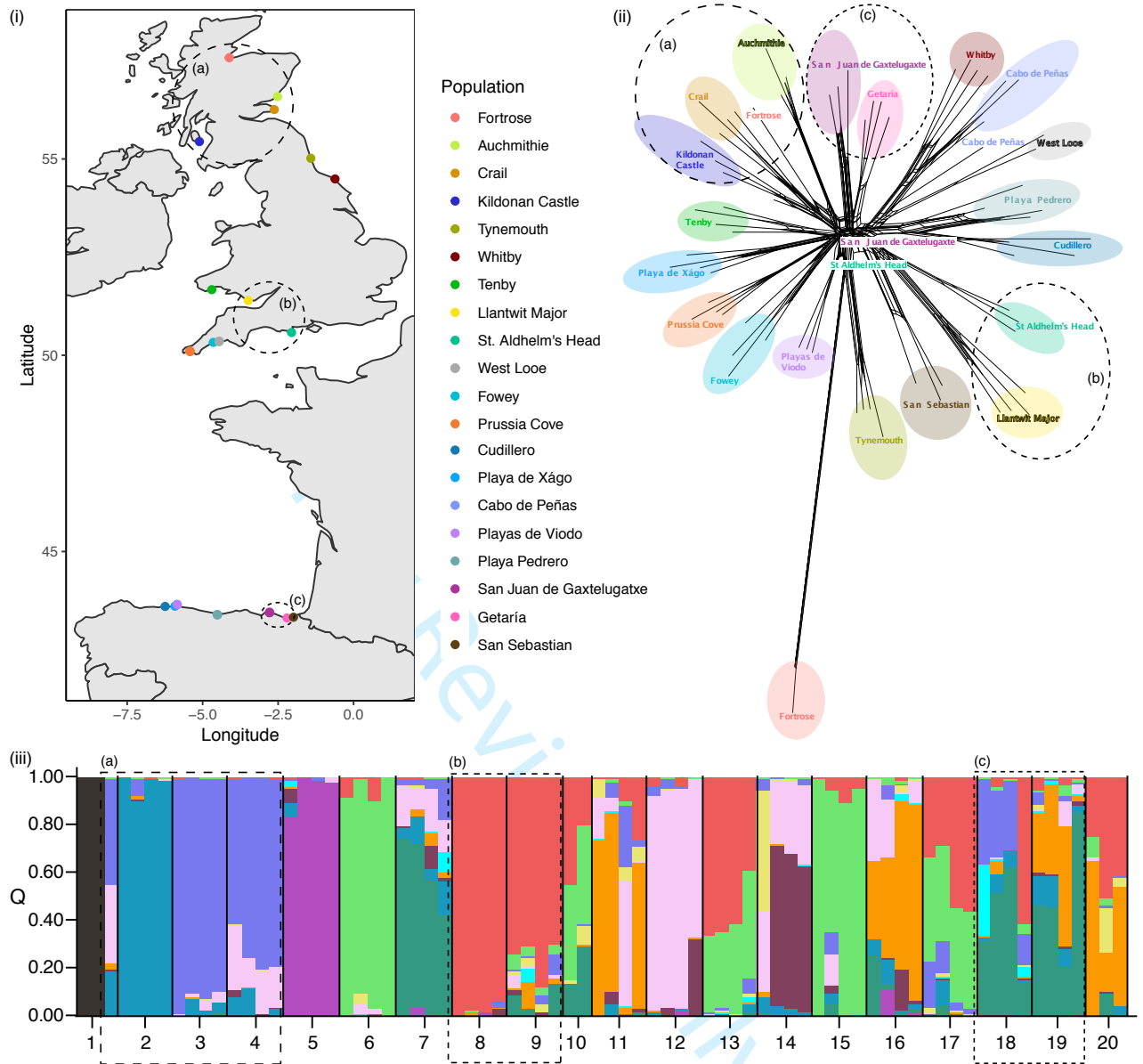


Figure 1: Population structuring of wild populations of *B. oleracea*. (i) Location of the populations considered here. (ii) Clustering of samples from RAXML (v8.2; GTRCAT model and 1000 maximum likelihood bootstrap replicates), visualised in SplitsTree4. (iii) STRUCTURE plot illustrating shared genetic ancestry for $K = 12$, ordered by population: 1 – Fortrose; 2 – Auchmithie; 3 – Craik; 4 – Kildonan Castle; 5 – Tynemouth; 6 – Whitby; 7 – Tenby; 8 – Llantwit Major; 9 – St. Aldhelm's Head; 10 – West Looe; 11 – Fowey; 12 – Prussia Cova; 13 – Cudillero; 14 – Playa de Xágo; 15 – Cabo de Peñas; 16 – Playas de Viodo; 17 – Playa Pedrero; 18 – San Juan de Gaxtelugatxe; 19 – Getaría; 20 – San Sebastian. Across the figures the same colours and numbering is used for each population. The dashed lines and letters indicate some clustering: (a) populations in Scotland; (b) populations closest to the Welsh-English border; and (c) populations in the Basque Country, Spain (excluding San Sebastian).

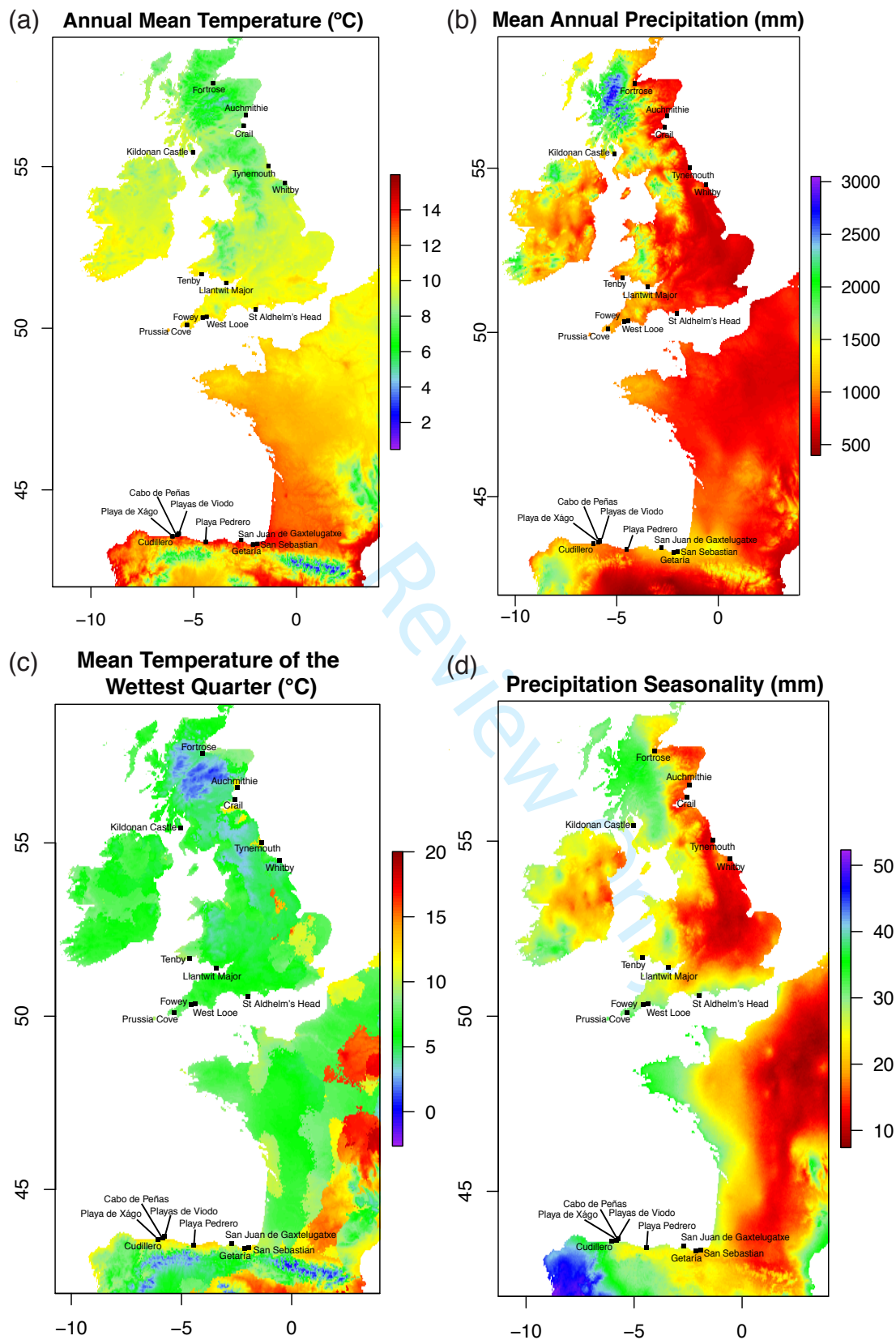
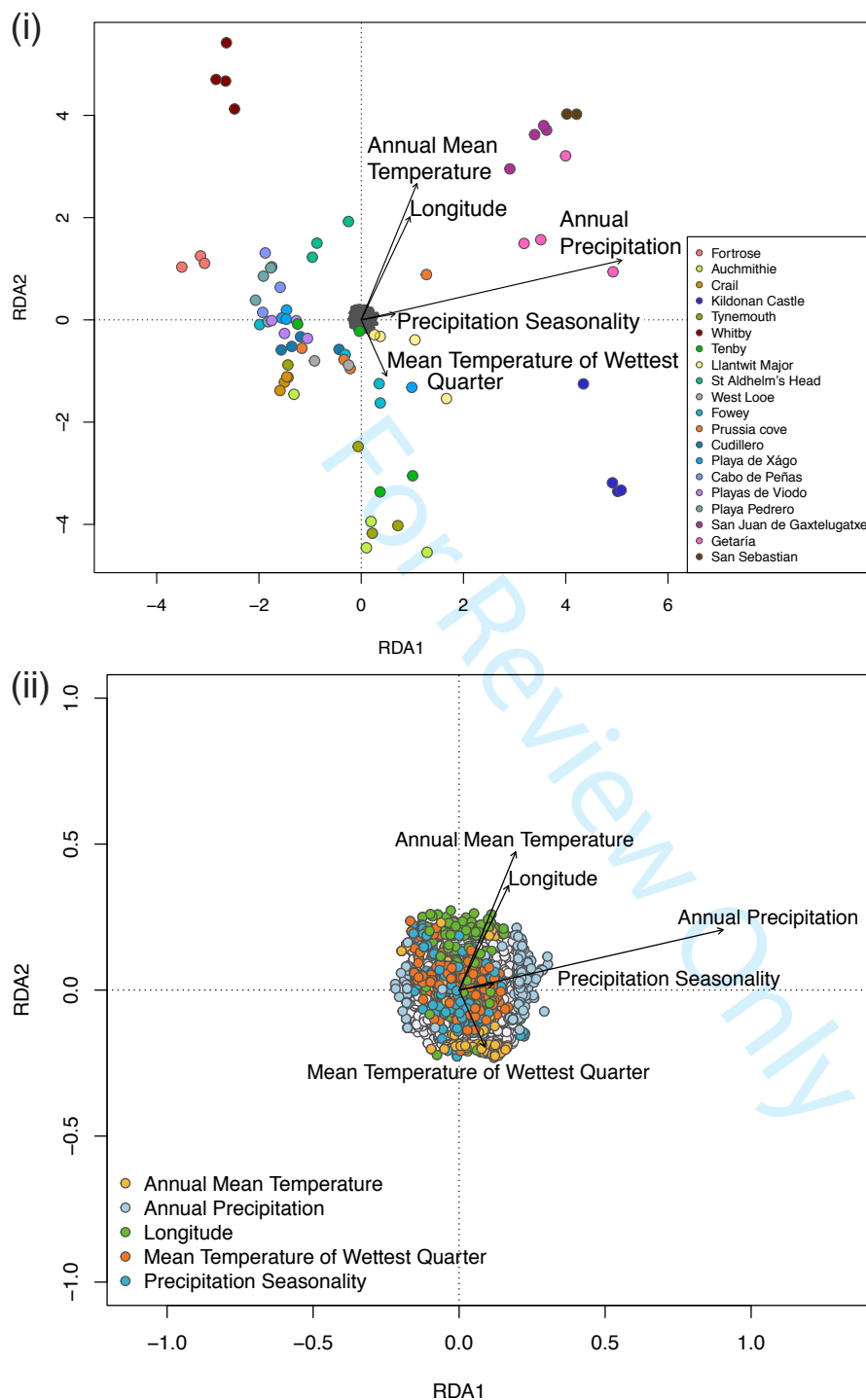


Figure 2: The distribution of sampled populations in relation to various climate variables: (a) annual mean temperature (°C); (b) mean annual precipitation (mm); (c) mean temperature of wettest quarter (°C); (d) precipitation seasonality (mm). These are averages between 1970 – 2000 obtained from the WorldClim database (Fick & Hijmans, 2017).

Table 3: Summary statistics of within *B. oleracea* population genetic diversity based on both variant nucleotide sites alone (var) and all sites (all) from dataset 1, showing: the number of individuals (N), the number of private alleles (PRI), expected heterozygosity (H_E), observed heterozygosity (H_O) and percentage of polymorphic loci (%).

Region†	Population	N	PRI	H_E		H_O		% All
				Var	All	Var	All	
ES	Authmithie	4	1683	0.1043	0.0012	0.1202	0.0014	33.0
ES	Crail	4	1727	0.1327	0.0019	0.1267	0.0018	52.8
ES	Fortrose	3	12951	0.2006	0.0032	0.1962	0.0031	76.4
WS	Kildonan Castle	4	1014	0.0903	0.0014	0.0944	0.0014	40.8
NEE	Tynemouth	4	1476	0.1023	0.0013	0.0881	0.0011	36.4
NEE	Whitby	4	1573	0.1200	0.0020	0.1184	0.0020	56.7
SW	Tenby	4	1568	0.1227	0.0014	0.1153	0.0013	40.5
SW	Llantwit Major	4	2073	0.1390	0.0023	0.1231	0.0022	66.2
SWE	Prussia Cove	4	1454	0.1019	0.0016	0.1064	0.0017	45.5
SWE	Fowey	4	1137	0.1126	0.0018	0.1083	0.0017	53.4
SWE	West Looe	2	1412	0.1150	0.0011	0.1328	0.0013	27.1
SWE	St. Aldhelm's Head	4	2470	0.1486	0.0014	0.1676	0.0016	39.4
A	Cudillero	4	716	0.0918	0.0015	0.0938	0.0016	44.3
A	Playa de Xágo	4	1583	0.1140	0.0012	0.1191	0.0012	33.4
A	Cabo de Peñas	4	698	0.0933	0.0015	0.0910	0.0014	42.5
A	Playas de Viodo	4	503	0.0545	0.0004	0.0580	0.0004	11.2
C	Playa Pedrero	4	1741	0.1313	0.0014	0.1408	0.0015	38.5
BC	San Juan de Gaxtelugatxe	4	2608	0.1423	0.0012	0.1471	0.0012	34.0
BC	Getaría	4	1550	0.1280	0.0021	0.1391	0.0023	59.8
BC	San Sebastian	3	2516	0.1530	0.0023	0.1538	0.0023	61.4

†Region codes: ES – East Scotland, WS – West Scotland, NEE – North-eastern England, SW – South Wales, SWE – South-western England, A – Asturias Spain, C – Cantabrica Spain, BC – Basque Country Spain.



53 **Figure 3:** (i) Redundancy analysis (RDA) ordination plot of the association between *B. oleracea* individuals (coloured points) and SNPs (dark grey points), with environmental variables. The different colours indicate which population each individual was from. (ii) RDA ordination plot of the SNPs alone, coloured for the environmental variable with which they were most strongly associated. For both (i) & (ii) the arrows indicate the environmental predictors and the strength of the association.

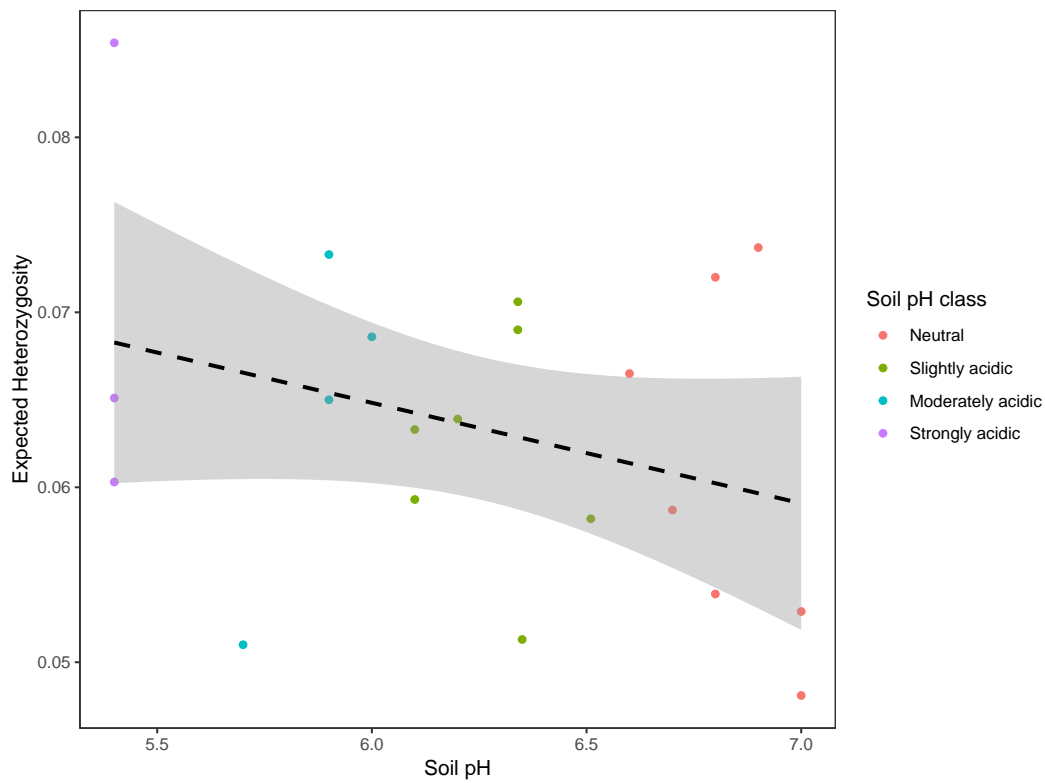


Figure 4: The relationship between expected heterozygosity and soil pH for 21 individuals from four soil pH classes categorised into: Neutral (6.6 - 7.3), Slightly acidic (6.1 - 6.5), Moderately acidic (5.6- 6.0) and Strongly acidic (5.0 - 5.5) based on USDA (1998). A linear model was used to fit a regression line (dashed black line), the standard error is shown in grey, p-value > 0.05.

Table 4: The top 18 candidate SNPs that mapped to unique genes in the *B. oleracea* reference genome and their annotations from 'Bolbase' (Yu et al., 2013).

Chromosome	Location	Identity	X	Bolbase gene name	Potential protein	Function
C09	32879582	1	-	Bol019890	Ribonucleotide reductase-related	Fatty acid metabolic process, creation of DNA from RNA
C04	39737611	0.999979	-	Bol021601	Unknown	
C09	8499546	1	+	Bol032146	Basic helix-loop-helix dimerisation region	Nucleus transcription regulation
C07	43014116	1	-	Bol042101	Toll-Interleukin receptor	Signal transduction, immune response, disease resistance
C02	233586	1	+	Bol012817	Laccase/multicopper oxidase	Copper ion binding, metabolic process, maybe formation and degradation of lignin
C04	22051514	0.999656	+	Bol044300	Protein kinase - serine/threonine	Protein kinase activity, signalling, plant defence
C03	29308196	0.472347	-	Bol012462	PIK-related kinase	Binding and DNA repair
C03	48963472	0.99438	+	Bol029900	Protein kinase	Protein kinase activity, signalling, plant defence
C04	28456859	0.999661	-	Bol009961	Cystathionine beta-synthase	Vitamin B6 pathway?
C03	9456274	1	-	Bol005573	Unknown	
C05	2317477	0.580051	-	Bol041075	Pentatricopeptide repeat	Often essential in mitochondria
C04	35972614	0.304057	+	Bol037830	Bacterial transferase haxapeptide repeat	Binding and transferase activity
C04	35104965	0.996501	+	Bol037950	Cyclin-like F-box	Growth and development
C03	2461137	0.999261	-	Bol034275	Serine/threonine-protein kinase	Signalling, plant defence
C02	233586	0.168963	-	Bol012816	Serine/threonine-protein kinase	Signalling, plant defence
C01	11164295	0.999978	+	Bol039465	Initiation factor eIF-4 gamma, MA3	
C01	11431159	1	+	Bol039505	Heat shock protein Hsp20	
C01	12106862	0.918256	-	Bol039585	F-box associated	

1
2
3
4 **Title: Feral populations of *Brassica oleracea* along Atlantic coasts in western Europe**
5
6
7

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

2
3 Running title: Feral *Brassica oleracea* in western Europe
4

5 Elizabeth A. Mittell^{1,2,*}, Christina A. Cobbold^{3,4}, Umer Zeeshan Ijaz⁵, Elizabeth A. Kilbride¹, Karen A.
6 Moore⁶ & Barbara K. Mable^{1,4}†

7
8 ¹Institute of Biodiversity, Animal Health and Comparative Medicine, University of Glasgow, UK; ²School
9 of Biology, University of St Andrews, UK; ³School of Mathematics and Statistics, University of Glasgow,
10 UK; ⁴The Boyd Orr Centre for Population and Ecosystem Health, University of Glasgow, Glasgow, UK;
11 ⁵School of Engineering, University of Glasgow, UK; ⁶Exeter Sequencing Service, University of Exeter,
12 UK

13
14 Corresponding authors: Elizabeth A. Mittell, *em294@st-andrews.ac.uk, e.mittell@gmail.com & Barbara
15 K. Mable, † Barbara.Mable@glasgow.ac.uk
16

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract

There has been growing emphasis on the role that crop wild relatives might play in supporting highly selected agriculturally valuable species in the face of climate change. In species that were domesticated many thousands of years ago, distinguishing wild populations from escaped feral forms can be challenging, but reintroducing variation from either source could supplement current cultivated forms. For economically important cabbages (Brassicaceae: *Brassica oleracea*), “wild” populations occur throughout Europe but little is known about their genetic variation or potential as resources for breeding more resilient crop varieties. The main aim of this study was to characterise the population structure of geographically isolated wild cabbage populations along the coasts of the UK and Spain, including the Atlantic range edges. Double-digest restriction-site associated DNA sequencing was used to sample individual cabbage genomes, assess the similarity of plants from 20 populations, and explore environment-genotype associations across varying climatic conditions. Interestingly, there were no indications of isolation-by-distance; several geographically close populations were genetically more

1
2
3
4
5 30 distinct from each other than to distant populations. Furthermore, several distant populations shared
6 31 genetic ancestry, which could indicate that they were established by escapees of similar source culti-
7
8 32 vars. **However**, there were signals of local adaptation to different environments, including a possible
9
10 33 relationship between genetic diversity and soil pH. Overall, these results highlight wild cabbages in the
11
12 34 Atlantic region as an important genetic resource worthy of further research **into their relationship with**
13
14 35 **existing crop varieties**.

15
16 36 **Keywords:** *Brassica oleracea*, feral populations, crop wild relatives, isolation-by-distance, environment-
17
18 37 genotype associations, domestication

21 22 38 **Introduction**

23
24
25 39 Domestication was an important transition within human societies, which allowed the rise of civilisations
26
27 40 (Diamond, 2002). Whilst vital for human success, there have been evolutionary consequences for the
28
29 41 domesticated organisms. In crop plants, the selection of 'domestication traits' has led to many desired
30
31 42 changes in physiological, morphological and life-history traits compared to their wild relatives (Milla,
32
33 43 Osborne, Turcotte, & Violle, 2015; Purugganan & Fuller, 2009). However, traits that are correlated with
34
35 44 those selected for (directly or indirectly) can also influence phenotypes via pleiotropic effects (Conner,
36
37 45 2002) and linkage disequilibrium (Falconer & Mackay, 1996). These genetic constraints and narrow
38
39 46 population bottlenecks can have unintended genetic consequences for crop plants, particularly elite
40
41 47 lines that are the result of intense artificial selection; e.g., reduced genetic diversity, increased genetic
42
43 48 drift and increased deleterious allele frequencies (Rauf, Teixeira da Silva, Khan, & Naveed, 2010; von
44
45 49 Wettberg et al., 2018). It is also likely that crop lines are constrained to some extent by the environment
46
47 50 within which they were originally domesticated. Therefore, to continue to utilise crop plants successfully,
48
49 51 it is important to understand both the genetic consequences of domestication, and where it occurred.

50
51 52 A classic example of domestication can be found in the commercially valuable species, *Brassica ol-*
52
53 53 *eracea* (recognised by Darwin, 1859; Walley et al., 2012). This single species contains a huge amount
54
55 54 of morphological diversity in cultivated varieties that has been around since at least the 1st Century (e.g.,
56
57 55 kale, kohlrabi, broccoli, Brussels sprouts and cauliflower; Maggioni, von Bothmer, Poulsen, & Lipman,
58
59 56 2018); the same morphological extremes are not found in wild populations. The origin of domesticated
60

1
2
3
4 57 *B. oleracea* crops and the 'wild' or 'feral' status of populations, found throughout the UK and along the
5
6 58 Atlantic coasts of north-western Europe (Raybould, Mogg, Clarke, Gliddon, & Gray, 1999), has been de-
7
8 59 bated in the literature (Allender, Allainguillaume, Lynn, & King, 2007; Gómez-Campo & Prakash, 1999;
9
10 60 Maggioni, 2015; Mitchell, 1976). Initially it was thought that different cultivars were independently do-
11
12 61 mesticated from wild populations on European Atlantic coasts (e.g., Spanish cabbage varieties were
13
14 62 domesticated from local wild Spanish populations; Gómez-Campo & Prakash, 1999), and that early
15
16 63 domesticates were introduced to and diversified within the Mediterranean region around 3,000 to 4,000
17
18 64 years ago (Allender et al., 2007). Information was limited when this hypothesis was favoured (Allender et
19
20 65 al., 2007; Gómez-Campo & Prakash, 1999), although there was already conflicting evidence (Mitchell,
21
22 66 1976). For example, Mitchell (1976) found that the locations of ancient human settlements and modern
23
24 67 *B. oleracea* populations coincided along UK coasts, providing a potential source of escapees from do-
25
26 68 mestic settings (agriculture or gardens) that could have established feral populations. This alternative
27
28 69 hypothesis that *B. oleracea* originated elsewhere and escaped into the wild in the Atlantic region has
29
30 70 been supported by recent linguistic and historical research (Maggioni, 2015; Maggioni et al., 2018).
31
32 71 Maggioni (2015) suggested that the most plausible hypothesis is that *B. oleracea* was domesticated in
33
34 72 the Mediterranean region, before being moved across Europe by people, where escaped plants estab-
35
36 73 lished now naturalised populations. However, the genetic status of *B. oleracea* in the Atlantic region is
37
38 74 still an open question (*B. oleracea* is classified as a native species in the UK and an alien species in
39
40 75 Spain; Euro+Med PlantBase, 2020).

41
42 76 The ease with which cultivated and wild *B. oleracea* plants can introgress is an issue for interpreting
43
44 77 variation within the *B. oleracea* species complex, as past hybridisation can obscure phylogeographic sig-
45
46 78 nals (Allender et al., 2007). However, for crop breeding purposes a close genetic relationship between
47
48 79 wild populations and domesticated cultivars may be seen as an advantage; higher genetic similarity
49
50 80 could make it easier to introgress adaptive traits from the wild into cultivated varieties (Hoisington et
51
52 81 al., 1999). An alternative view is that if these populations are feral they would have experienced the
53
54 82 same domestication bottleneck as many cultivars (von Wettberg et al., 2018), and therefore they may
55
56 83 not be the important reservoirs of genetic diversity that crop wild relatives are typically assumed to be.
57
58 84 Compared to domestication, feralization is under-investigated; however modern genomic data are al-

1
2
3
4 85 lowering its occurrence to be identified and consequences better understood (see examples in Henriksen,
5
6 86 Gering, & Wright, 2018). Despite the agricultural importance of *B. oleracea*, there has not yet been a
7
8 87 comprehensive genetic analysis of wild populations in the Atlantic region that would allow assessment
9
10 88 of their utility as sources of variation for cultivation.

11
12 89 Escaped plants can be thought of as 'invasive' species, which are defined as those that became
13
14 90 established after introduction outside of the biogeographic region within which they evolved (Prentis,
15
16 91 Wilson, Dormontt, Richardson, & Lowe, 2008). However, it is not always clear where these 'native'
17
18 92 regions are located, as is the case of *B. oleracea*, or why certain species are successful where others
19
20 93 are not. Furthermore, wild populations of *B. oleracea* do not have the characteristics that are thought to
21
22 94 be important for successful establishment in novel locations (i.e. 'invasive traits'; Funk, Standish, Stock,
23
24 95 & Valladares, 2016). For example, wild *B. oleracea* are: perennials rather than annuals, woody rather
25
26 96 than herbaceous, relatively slow- rather than fast-growing, and predominantly outcrossing rather than
27
28 97 selfing. Self-fertilisation in plants is inhibited by polymorphic self-incompatibility (SI) recognition systems
29
30 98 where haplotype blocks encode distinct proteins for pollen-pistil recognition (Charlesworth, Vekemans,
31
32 99 Castric, & Glémin, 2005). A strong SI system exists in *B. oleracea* (a single-locus system with over
33
34 100 60 alleles; Raybould et al., 1999), making them predominantly self-incompatible (Kitashiba & Nasrallah,
35
36 101 2014; Walley et al., 2012; Yousef, Mueller, Börner, & Schmid, 2018). Development of self-compatible
37
38 102 lines can aid in propagation of cultivated forms (e.g., Xiao et al., 2019), but reduce adaptive potential
39
40 103 to changing environmental conditions. Therefore, even if the "wild" populations include escaped forms,
41
42 104 retention of a wide range of self-incompatibility alleles could be used to enhance the potential of breeding
43
44 105 strategies designed to maintain heterosis.

45
46 106 Currently too little is known about levels of genetic variation and population structure in wild *B. ol-*
47
48 107 *eracea* populations to fully assess the potential for use of plants from different regions to supplement
49
50 108 crop diversity. Population structure and within population genetic diversity are impacted by gene flow,
51
52 109 which occurs via pollen and seeds in plants (Scheepens, Frei, Armbruster, & Stöcklin, 2012; Slatkin,
53
54 110 1987). The main pollinators of *B. oleracea* are bees that fly short distances between plants (average
55
56 111 2 m; Raybould et al., 1999). Seed dispersal was previously thought to be limited to approximately 4 m
57
58 112 (Watson-Jones, Maxted, & Ford-Lloyd, 2006). However, Wichmann et al. (2009) found that wind can

1
2
3
4 113 spread seeds up to 250 m, and that rare-long distance dispersal events of up to 10 km could occur
5
6 114 if seeds became attached to people's shoes. Therefore, although gene flow may be limited between
7
8 115 geographically close populations leading to high genetic structuring in some instances, in other cases,
9
10 116 such as where plants grow close (0 - 4 m) to well used coastal paths, gene flow might be greater than
11
12 117 expected. Genetic diversity estimates have been made in some *B. oleracea* populations within the
13
14 118 Atlantic region (e.g., Table 1), but the northern edge (Scotland) has not been investigated. A correla-
15
16 119 tion between genetic distance and geographic distance in wild *B. oleracea* populations was found in
17
18 120 some studies (Raybould et al., 1999; Sánchez-Yélamo, 2014) but not others (Christensen et al., 2011;
19
20 121 Watson-Jones et al., 2006). Interestingly, Watson-Jones et al. (2006) also considered some environ-
21
22 122 mental variables and found that higher soil pH was associated with lower genetic diversity in English
23
24 123 and Welsh populations. The inconsistency in previous studies could be due to the varying spatial scales
25
26 124 and molecular markers used. However, overall, these results highlight the uncertainty in the status and
27
28 125 genetic contents of wild *B. oleracea* populations in the Atlantic region, as well as the potential effect of
29
30 126 environment on the plant genetics. Filling these knowledge gaps could provide important insights into
31
32 127 these crop wild relatives for agricultural use.

33
34 128 *Brassica oleracea* is a good model for investigating the genetic resources available (e.g., the extent
35
36 129 of genetic diversity and local adaptation) in a potentially feral crop wild relative because it is diploid
37
38 130 and a reference genome is available (Liu et al., 2014). Therefore, compared to other crop species
39
40 131 (e.g., polyploids) genetic analyses are simpler. For many questions whole-genome sequencing is un-
41
42 132 necessary (Rockman, 2012) and reduced-representation methods, such as double-digest restriction
43
44 133 associated DNA sequencing (ddRADseq), are sufficient to: assess genetic diversity within and between
45
46 134 populations (Andrews, Good, Miller, Luikart, & Hohenlohe, 2016); determine population genetic struc-
47
48 135 turing (Gao et al., 2017); and investigate potential associations between genotypes and environmental
49
50 136 variables (Forester, Lasky, Wagner, & Urban, 2018). Therefore, ddRADseq is an appropriate method
51
52 137 for considering the genetic resources in, and local adaptation of, *B. oleracea* populations across their
53
54 138 Atlantic range.

55
56 139 Overall, current knowledge on genetic variation of *B. oleracea* in wild populations is patchy in geo-
57
58 140 graphic coverage and based on outdated molecular genetic techniques (Table 1). Therefore, this study
59
60

1
2
3
4 141 combined modern genetic techniques and the reference genome available for this species to increase
5
6 142 the power to detect differences among populations across a broad geographic range. The following
7
8 143 questions were addressed: (1) how much genetic variation exists among wild populations of *B. oler-*
9
10 144 *acea* in the UK and Spain; (2) how are populations structured in the Atlantic region and how much
11
12 145 differentiation exists between isolated populations; and (3) are there signals of local adaptation to the
13
14 146 environment? The results provide insights into the utility of *B. oleracea* as a crop wild relative genetic
15
16 147 resource for agriculture, as well as shed light on the most likely region of *B. oleracea* domestication.
17
18
19

20 148 **Materials and Methods**

21
22
23
24 149 Twenty-four populations of *B. oleracea* were chosen from the UK and Spain to cover both a latitudi-
25
26 150 nal and longitudinal gradient of the Atlantic range for genetic analyses (Figure 1i & Table 2). French
27
28 151 populations were not sampled here, but are the focus of a recent genetic analysis by Maggioni *et al.*
29
30 152 (personal communication). Leaves were collected from four individual plants from each population for
31
32 153 DNA extraction, as has been successfully applied to the study of population structure in wild relatives
33
34 154 in the Brassicaceae (Buckley, Holub, Koch, Vergeer, & Mable, 2018). Nazareno, Bemmels, Dick, and
35
36 155 Lohmann (2017) found that compared to “traditional” population genetic markers these smaller sample
37
38 156 sizes are sufficient for various population statistics when large numbers of SNPs are available. The
39
40 157 bedrock for each population was obtained from the British Geological Survey (BGS, 2018) and the Insti-
41
42 158 tuto Geológico y Minero de España (IGME, 2018). The first year a written record of a population exists
43
44 159 was obtained for the UK populations from the Botanical Society of Britain & Ireland (BSBI, 2018). No
45
46 160 equivalent records could be found for the Spanish populations.
47
48

49 161 **Molecular methods**

50
51
52 162 High molecular weight DNA was extracted from the leaves of 96 individuals from 24 populations (Table
53
54 163 2) using DNeasy Plant Mini Kits (QIAGEN, Hilden, Germany) and quantified using a Qubit 2.0 Fluorom-
55
56 164 eter (ThermoFisher Scientific, Waltham, Massachusetts, U.S). Four samples from each population were
57
58 165 sent for library preparation and sequencing at University of Exeter Sequencing Service. Double-digest
59
60

1
2
3
4 166 RADseq libraries were made using a modification of the method in Wu et al. (2016) that allowed Nex-
5
6 167 teraXT indexes (Illumine Corp., USA) to be used for multiplexing samples. In addition, an RYRY spacer
7
8 168 was inserted in the adapter 3' of the Illumina sequencing primer annealing site to provide additional
9
10 169 complexity at the start of read 1 immediately before the Sac1 sticky end. For each sample 400 ng DNA
11
12 170 was fully digested with Sac1 and Mse1 restriction endonucleases and purified using Ampure XP beads.
13
14 171 Illumina compatible i5 adapters were designed to ligate to the at the AGCT-3' sticky end left after Sac1
15
16 172 digest, and Illumina compatible i7 adapters were designed to ligate to the 5'-TA overhangs remaining
17
18 173 after Mse1 digest. Adapter-ligation excess adapters were removed using Ampure XP beads. DNA frag-
19
20 174 ments were amplified by 12 cycles of indexing PCR, purified, size selected (inserts 330-670 bp) and
21
22 175 validated using a Tapestation D1000 HS Screentape (Agilent Technologies Ltd). Libraries were equimo-
23
24 176 lar pooled and the pool concentration was calculated after qPCR. Libraries were denatured, diluted and
25
26 177 sequenced with 125bp paired-end reads on Illumina HiSeq 2500 using SBS High Output reagents v4
27
28 178 (Illumina Corp., USA).

179 Data processing

180 Reads were demultiplexed and trimmed to 100 bp using cutadapt (Martin, 2011). These were then
181
182 cleaned and quality filtered using the *process_radtags* pipeline in Stacks v1.47 (Rochette & Catchen,
183
184 2017). Bowtie (v2; Langmead & Salzberg, 2012) and samtools (v1.9; Li et al., 2009) were used to
185
186 align the reads to the *B. oleracea* reference genome (Liu et al., 2014). A catalogue of stacks was then
187
188 created using *ref_map* (Stacks) with the default settings. The *populations* pipeline (Stacks) was used to
189
190 filter the data, and calculate summary statistics. Three datasets were generated with different filtering
191
192 parameters depending on the downstream analysis. Firstly, for dataset 1 (*within individuals*), which was
used to estimate genetic diversity within individuals and in phylogenetic analyses, all individuals were
filtered as a single population, and loci were retained if they had a minimum individual stack depth of
five, a minimum minor allele frequency of 0.01, a maximum observed heterozygosity of 0.7 and were
present in 60% of individuals. Secondly, dataset 2 was generated using the same filtering as dataset 1
but SNPs linked within each RAD locus were avoided by only retaining one SNP at random per locus;
required for population structure analyses (Pritchard, Stephens, & Donnelly, 2000). Finally, for dataset

1
2
3
4 193 3 (*within populations*), which was used to calculate genetic distance between populations, individuals
5
6 194 were assigned to their population of origin and loci were retained if present in 50% of the populations.
7
8 195 This filtering was designed to reduce the inclusion of duplicate loci and balance the amount of missing
9
10 196 data with the number of informative loci (Andrews et al., 2016). A minimum stack depth of five is higher
11
12 197 than the default of two, but within the recommended range (Paris, Stevens, & Catchen, 2017), and
13
14 198 helps to remove potential paralogues. Spurious SNPs were avoided by using a minor allele frequency
15
16 199 of > 0.01 (Marandel et al., 2020), and the combination of a maximum observed heterozygosity of 0.7
17
18 200 (70% of the individuals or populations can be heterozygous for each locus) which are present in either
19
20 201 60% of individuals (datasets 1 and 2) or 50% of the populations (dataset 3) retains loci that have been
21
22 202 successfully genotyped across individuals, but are not completely heterozygous. The summary statistics
23
24 203 for each population were calculated in Stacks during the filtering of dataset 3 and included: the number
25
26 204 of private alleles (PRI), expected heterozygosity (H_E), observed heterozygosity (H_O), percentage of
27
28 205 polymorphic loci (%; Table 3), the inbreeding coefficient (F_{IS}) and nucleotide diversity (π ; Supplementary
29
30 206 information).

207 Data analyses

31
32
33
34
35
36 208 Clustering of samples within and between populations was investigated with dataset 1 using RAxML
37
38 209 (v8.2; GTRCAT model and 1000 maximum likelihood bootstrap replicates; Stamatakis, 2014) and visu-
39
40 210 alisation in SplitsTree4 (Huson & Bryant, 2005). To estimate the number of putative genetic clusters
41
42 211 (K) and assess shared genetic ancestry, STRUCTURE (v2.3.4; Pritchard et al., 2000) was used with
43
44 212 dataset 2, so as not to inflate sharing based on multiple SNPs within a RAD locus. A range of K values
45
46 213 were tested (the number of populations successfully sequenced plus one; 1 – 21) using an admixture
47
48 214 model that assumed correlated allele frequencies. For each K, five independent replicates of 100,000
49
50 215 MCMC repetitions, after a burn-in of 10,000 iterations, were run. The most likely K was selected us-
51
52 216 ing the log likelihoods and deltaK (Evanno, Regnaut, & Goudet, 2005). To see if there were significant
53
54 217 differences between estimates of H_E and H_O , pairwise-ANOVAs were carried out in R version 3.4.0 (R
55
56 218 Core Team, 2017) on estimates from dataset 3 based on variant sites alone and all sites. A genetic
57
58 219 distance matrix was created using dataset 3, and the latitude and longitude of each population was
59
60

1
2
3
4 220 used to calculate a geographic distance matrix using 'Haversine' Great Circle Distance in the R package
5
6 221 'geosphere' (Hijmans, 2017). In addition, genetic and geographic matrices were created for Spanish
7
8 222 and UK populations **separately**, alongside a temporal distance matrix for the year when each population
9
10 223 was first recorded within the UK (first population record; Table 2). Mantel tests were carried out with
11
12 224 9999 replicates on the region-wide matrices and country matrices separately, to assess both the overall
13
14 225 and within country isolation-by-distance. Mantel tests were also carried out on the UK specific matri-
15
16 226 ces to investigate any relationship between the first population records and the genetic and geographic
17
18 227 distances.

20 228 A subset of dataset 1 where the soil pH was known was used to investigate the relationship between
21
22 229 soil pH and H_E – e.g., is a higher soil pH associated with lower genetic diversity? A linear model with
23
24 230 soil pH as a predictor variable and H_E as a response variable was run on 21 individuals (**across six**
25
26 231 **populations**) from four soil pH classes: Neutral (6.6 - 7.3), Slightly acidic (6.1 - 6.5), Moderately acidic
27
28 232 (5.6 - 6.0) and Strongly acidic (5.0 - 5.5) based on USDA (1998).

30 233 In order to identify potential genotype-environment associations, redundancy analyses (RDA) were
31
32 234 carried out using dataset 1 following Forester et al. (2018) with the R packages 'vegan' and 'pysch'
33
34 235 (Oksanen et al., 2017; Reville, 2018). The climate dataset was downloaded from the WorldClim
35
36 236 database at a resolution of 4.5 km (Fick & Hijmans, 2017). This dataset is based on measurements
37
38 237 made between 1970 – 2000. Therefore, it is assumed that any changes in climate will be consistent
39
40 238 enough across the study gradient to maintain differences in the averages and variation between pop-
41
42 239 ulations. The 19 climate variables available from WorldClim for our dataset were checked for pairwise
43
44 240 correlations and the estimated variance inflation factor (VIF). Variables with correlations $> |0.7|$ and
45
46 241 $VIF > 10$ were removed, leaving: 'Annual Mean Temperature', 'Mean Temperature of Wettest Quarter',
47
48 242 'Annual Precipitation' and 'Precipitation Seasonality'. Longitude was included as an additional predictor
49
50 243 variable because it was weakly correlated with climatic variables. Those SNPs that had RDA load-
51
52 244 ings with q-values < 0.1 were considered outlier loci, and were compared to the annotated *B. oleracea*
53
54 245 genome using Bedtools (v.2.17.0; Quinlan & Hall, 2010), followed by a search of the online resource
55
56 246 'Bolbase' (Yu et al., 2013) to investigate putative gene functions.

Results

Patterns of genetic diversity

A total of 115,746,909 reads from 76 individuals (20 populations; Table 2) were of sufficient quality and retained for down-stream analysis (average reads per individual: 1,522,986; range: 220,363 – 5,361,799; Supplementary Table 1). For four of the populations, no individuals were successfully sequenced and so these were not included in these analyses. On average 86.3% (range 82.5 - 88.6) of reads mapped to the reference genome (Supplementary Figure 1). Datasets 1 and 2 contained 42,517 and 13,352 SNPs, respectively, across 13,352 RAD-loci (Supplementary Table 2). There were 140,131 SNPs across 53,539 RAD-loci in dataset 3 (Supplementary Information).

Based on variable nucleotide sites only (Table 3), average estimates of genetic diversity (considering H_E) were lower than in the studies cited in Table 1; the average across populations was 0.120 among both UK (range 0.090 – 0.200) and Spanish (range 0.055 – 0.153) populations. Observed heterozygosity was consistently significantly (H_O $p < 0.001$) greater than H_E for all populations and average F_{IS} was similar in the two geographic regions (UK: average = 0.039, range = 0.001 to 0.084; Spain: average = 0.027, range = 0.025 to 0.031). There was thus no evidence of inbreeding (as expected given the genetically controlled self-incompatibility system) but heterozygosity excess was apparent in all populations. The Fortrose population contained 10-fold more private alleles compared to all other populations and had the highest values for both H_E and H_O . Values considering all sites were lower but did not change conclusions about relative patterns of diversity (Table 3).

Population structure

Based on the RAxML tree, the majority of individuals clustered by population, with the exceptions of: (i) two individuals that did not cluster with any population (one in San Juan de Gaxtelugatxe, Spain and one in St Aldehelm's Head, UK), and (ii) an individual from Fortrose (Scotland, UK) that clustered more closely with other Scottish populations than other individuals from Fortrose (Figure 1ii). The most likely number of genetic clusters from STRUCTURE analyses was $K = 12$. Most individuals were admixed, however, six of the UK populations (Fortrose, Auchmithie, Crail, Tynemouth, Whitby and Llantwit Major)

1
2
3
4 273 were dominated by a single genetic ancestry, and two individuals from Fortrose were distinct from both
5
6 274 the third individual from Fortrose and all other samples (Figure 1iii). The dominant genetic ancestry
7
8 275 seen in individuals from Whitby (UK) also dominated the genetic ancestry of individuals from Cabo de
9
10 276 Peñas (Spain), and similarly, the dominant genetic ancestry seen in individuals from Tenby (UK) was
11
12 277 most prevalent in individuals from San Juan de Gaxtelugatxe and Getaría (Spain). There were three
13
14 278 potential regional clusters indicated by the RAxML tree and STRUCTURE analysis: (a) populations
15
16 279 in Scotland; (b) populations closest to the Welsh-English border; and (c) populations in the Basque
17
18 280 Country, Spain (excluding San Sebastian; Figure 1). However, the clustering of populations was not well
19
20 281 resolved and these 'regional clusters' were not always the most geographically close populations (e.g.,
21
22 282 in cluster c, San Sebastian is closer to Getaría than San Juan de Gaxtelugatxe geographically but not
23
24 283 genetically). No isolation-by-distance was predicted by the data either region-wide, or within Spanish
25
26 284 or UK populations alone (Mantel test p-values = 0.474, 0.658 and 0.705, respectively). Furthermore,
27
28 285 no relationship was found between the first record for each of UK populations (Table 2) with either
29
30 286 geographic or genetic distance (Mantel test p-values = 0.114 and 0.933, respectively).
31
32
33

34 287 **Environmental associations**

35
36 288 Overall, environmental variables explained 2.3% (adjusted r-squared) of the variation in the SNPs using
37
38 289 RDA analysis; the strongest association of genotype with the environment was with annual precipitation
39
40 290 (Figure 3). This environmental variation was strong enough to be reflected in the clustering of individ-
41
42 291 uals, including the genetically distinct individuals from Fortrose (UK; Figure 1iii). For example, across
43
44 292 regions, west Scotland and the Basque country experienced the greatest amount of annual precipitation
45
46 293 on average (Figure 2b), whereas the annual mean temperature was greater in the Basque country com-
47
48 294 pared to west Scotland (Figure 2a). Individuals from populations in these regions separated from other
49
50 295 populations in the same direction as annual precipitation, but in opposing directions in relation to annual
51
52 296 mean temperature (Figure 3i). Individuals from Whitby (UK) appear to have experienced a colder, drier
53
54 297 environment than the geographically closest population, Tynemouth (UK), which was also reflected in
55
56 298 the RDA analysis. Linear modelling indicated a non-significant negative trend between genetic diversity
57
58 299 (H_E , H_O & π) and soil pH (i.e. plant genetic diversity decreased as soil pH increased. Only H_E is shown
60

1
2
3
4 but the same relationship was found with H_O & π ; Figure 4).
5

6 There were 2249 unique candidate SNPs associated with the predictor variables from the RDA anal-
7
8 ysis; the majority of these (1039) were most closely associated with 'Mean Temperature of Wettest
9
10 Quarter', followed by 'Precipitation Seasonality' (349), 'Longitude' (333), 'Annual Precipitation' (269) and
11
12 'Annual Mean Temperature' (259). These were fairly evenly distributed across the genome with no indi-
13
14 cations of any single SNP with a large effect. A few SNPs that were more closely associated with annual
15
16 precipitation had strong loadings along axis 1 in the direction of the annual precipitation vector (Figure
17
18 3(ii)). In total, 221 candidate SNPs mapped to unique genes in the *B. oleracea* reference genome, and
19
20 of the top 18, six were annotated as part of the receptor-like kinase family (Table 4).
21
22
23

24 Discussion 25

26
27 The results presented here provide the first genome-wide estimates of genetic variation and population
28
29 genetic structure of wild cabbages collected from across the UK and Spain. Although direct compar-
30
31 isons with cultivated species would be required to rigorously test hypotheses about origins of these
32
33 populations, patterns of variation are consistent with recent linguistic and historical evidence (Maggioni,
34
35 2015; Maggioni et al., 2018) suggesting that the domestication of *B. oleracea* crops occurred in the
36
37 Mediterranean, domesticates were moved by people across Europe, escaped and established wild pop-
38
39 ulations in the Atlantic region. For example, there was no indication of isolation-by-distance from north-
40
41 ern Scotland to Spain ($> 14^\circ$ latitude), which might be expected if these plants were natural colonisers
42
43 following common phylogeographic patterns (e.g., Sharbel, Haubold, & Mitchell-Olds, 2000). Further-
44
45 more, genetic ancestry and clustering analyses suggested that geographically distant populations may
46
47 have similar genetic sources, and could therefore have been established by similar source cultivars.
48
49 The consistent excess of heterozygotes across populations, combined with evidence for admixture from
50
51 STRUCTURE analyses, suggests mixing between 'isolated' populations (Rousset & Raymond, 1995),
52
53 which could be due to interbreeding between cultivated plants growing near the wild populations. This
54
55 highlights the possibility of continued introgression between cultivated and wild plants. Despite the lack
56
57 of geographic genetic population structuring, there were signals of local adaptation to different climates
58
59
60

1
2
3
4 326 based on RDA analyses. In addition, within population genetic diversity estimates were comparable to
5
6 327 other studies (e.g., Christensen et al., 2011; Watson-Jones et al., 2006), and as Watson-Jones et al.
7
8 328 (2006) found, lower genetic diversity estimates were associated with higher soil pH. Therefore, these
9
10 329 wild populations could hold useful adaptive alleles for plant breeding, and a suitable approach to investi-
11
12 330 gate traits of agricultural interest (e.g., drought tolerance) could be to choose populations based on their
13
14 331 environment of origin. However, further sequencing of a range of cultivars from different geographic
15
16 332 regions would be required to further test these hypotheses.

333 **Patterns of Genetic Diversity**

334 Although the magnitude of estimates of genetic diversity based on the ddRADseq data presented here
335 were lower than in previous studies (see Table 1) using allozymes (Lanner-Herrera, Gustafeson, Filt, &
336 Bryngelsson, 1996; Lázaro & Aguinagalde, 1998; Raybould et al., 1999), microsatellites (Raybould et
337 al., 1999) or AFLPs (Watson-Jones et al., 2006; Christensen et al., 2011), patterns of variation within
338 the UK and Spain were strikingly similar to one another. Most populations also showed a relatively con-
339 sistent excess of heterozygosity. These similarities could provide evidence for relatively recent origins
340 of populations in the two regions, but whether this was from feralisation of cultivars or natural differ-
341 entiation after natural colonisation cannot be distinguished by the data. Although there has been an
342 ongoing debate as to the origin of wild *B. oleracea* populations in the Atlantic region (Song, Osborn,
343 & Williams, 1990; Allender et al., 2007; Maggioni, 2015), domestication of *B. oleracea* in the Mediter-
344 ranean region has been suggested by other genetic, phenotypic and linguistic studies (Mitchell, 1976;
345 Maggioni, 2015; Maggioni et al., 2018). The subsequent movement of *B. oleracea* cultivars across Eu-
346 rope could then have resulted in a much narrower bottleneck than the initial domestication bottleneck
347 in the Mediterranean as it removed the chance of gene flow from the wild relatives they originated from
348 (Kofsky, Zhang, & Song, 2018). Consistent with this hypothesis, although the putative Mediterranean
349 progenitor species remains unknown, Allender et al. (2007) found much greater estimates of genetic
350 diversity within potential progenitor species from the Mediterranean region than either previous genetic
351 diversity estimates made in *B. oleracea* (e.g., Christensen et al., 2011; Watson-Jones et al., 2006) or in
352 this study.

Population structure

Several of the analyses here suggest less population structuring than might be expected in such geographically distinct populations if natural range expansion followed by isolation occurred. In this dataset, since the first recorded population (Tenby in 1773), one to three new populations have been recorded every thirty years within the UK (Table 2). However, neither the date the UK populations were first recorded, nor the genetic distances between populations in the UK and Spain, had a geographical pattern (i.e. no isolation by distance). Furthermore, although the majority of individuals clustered by population and some regional clustering was seen (Figure 1), it would not be possible to predict whether two individuals from geographically close or geographically distant populations are more genetically similar to each other. For example, Fowey and Prussia Cove (UK populations), and West Looe and Cabo de Peñas (UK and Spanish populations respectively), clustered together and shared more genetic ancestry than Fowey and West Looe, which are the closest geographically. Although more sampling would be required to explicitly test it, the evidence here suggests that these plants have not colonised the Atlantic region following common phylogeographic patterns (e.g., Sharbel et al., 2000) and therefore is consistent with *B. oleracea* domestication occurring outside of the Atlantic region. This is in line with results from other genetic, phenotypic and linguistic studies, which suggest the Mediterranean region is the most likely location for *B. oleracea* domestication (Maggioni, 2015; Maggioni et al., 2018; Mitchell, 1976).

The genetic ancestry and clustering analyses hint that populations could have been established by escapees from different cultivars. The majority of individuals were assigned to multiple sources of genetic ancestry (Figure 1iii), however, there were also cases where one putative source dominated at the individual- and population-levels, which could be the overall genetic background from the original source cultivar. Interestingly, there were two distinct individuals from Fortrose (10-fold more private alleles than other populations; Table 3) with a source that was assigned to no other individuals. Due to the ease of interbreeding between cultivars (Allender et al., 2007), this could indicate that these two Fortrose individuals are recent escapees from a different source population (e.g., local gardens), which are yet to have mixed with other individuals within the population. Furthermore, the more recent record of the population at Fortrose (1968), and the lack of assignment to other populations, suggests that this genetic background could be from a cultivar that has not been grown for a long period of time or widely

1
2
3
4 381 around the Atlantic coastlines. The excess of heterozygotes (H_O was significantly greater than H_E) and
5
6 382 the general mix of shared genetic ancestry across such a wide geographical area in distinct populations,
7
8 383 could also be an indication of continued introgression into these wild populations from agricultural and
9
10 384 horticultural sources. It would be interesting to identify popular cultivars in the local areas of these
11
12 385 populations, including any changes in the preferred cultivars through time, to investigate patterns of
13
14 386 introgression in more detail. Such direct comparisons with cultivars could identify the most likely founder
15
16 387 of these populations.

17
18 388 Using chloroplast microsatellite DNA markers, Allender et al. (2007) found two haplotypes in *B.*
19
20 389 *oleracea* around the coasts of the UK; out of sixteen populations, fourteen were C:01 and two were
21
22 390 C:04. The two populations with the C:04 haplotype were in Tyne & Wear, in the northeast of England;
23
24 391 in the current study, this area is represented by the Tynemouth and Whitby populations. In line with the
25
26 392 rarity of the chloroplast haplotypes identified in this region in the previous study, these two populations
27
28 393 clustered most closely with populations not sampled by Allender et al. (2007); Tynemouth clustered with
29
30 394 Fortrose, Scotland, and Whitby with the Spanish population Cabo de Peñas. Based on this information,
31
32 395 it might be expected that the chloroplast haplotypes of Fortrose and Cabo de Peñas would also be
33
34 396 C:04. In addition, the C:01 haplotype found in the majority of the UK populations was also found in
35
36 397 four other species of *Brassica* (Allender et al., 2007), suggesting either that this is the ancestral form or
37
38 398 introgression between species. A combination of nuclear and chloroplast information could be useful for
39
40 399 disentangling the population histories further, particularly in relation to identifying introgression.

41
42 400 Knowledge of the founding cultivars would be useful for both plant breeders and those interested
43
44 401 in invasive species. It could provide insights into how different cultivars have adapted (and therefore
45
46 402 may adapt in the future) to different environmental conditions, and could also be thought of as a way
47
48 403 to compare invasion success within a species. *Brassica oleracea* lack the characteristics thought to be
49
50 404 fundamental for establishment in novel locations (invasions; Funk et al., 2016), but perhaps amongst
51
52 405 the huge phenotypic variation found within this species, some traits are more likely to lead to successful
53
54 406 'invasions' of particular cultivars compared to others. For example, a cultivated Danish kale was the
55
56 407 most likely source for a wild population found in Denmark (based on AFLP markers; Christensen et
57
58 408 al., 2011), and it could be that all the Atlantic populations were established by different kale cultivars.

1
2
3
4 Overall, populations of *B. oleracea* growing along Atlantic coasts would be an excellent study system to
5
6 improve understanding of invasive species that are likely to harbour useful adaptive traits for agriculture.
7

8 While comparisons with published whole genome sequence data or other types of genotype by
9
10 sequencing approaches (e.g., Stansell et al., 2018) for cultivated *B. oleracea* would be interesting to
11
12 more explicitly test origins of the populations studied here, there are several issues with ddRAD data
13
14 that would make this challenging and potentially hard to interpret. A benefit of ddRAD sequencing is
15
16 the generation of discrete loci that are standardised to the same length. However, the resulting short
17
18 sequence segments normally contain only one or a few SNPs, which does not allow accurate assign-
19
20 ment of paralogs in highly duplicated and rearranged genomes such as found in the Brassicaceae (e.g.,
21
22 Schranz, Lysak, & Mitchell-Olds, 2006). Instead, filtering pipelines to allow population genetics analyses
23
24 based on ddRAD data are designed to be conservative (Paris et al., 2017; Marandel et al., 2020). This
25
26 filtering results in fewer loci retained, but it should reduce risks of including duplicates. In the current
27
28 study, excess heterozygosity was observed consistently across populations, which could suggest his-
29
30 torical introgression. Although we cannot completely rule out the influence of combining duplicates (Ilut,
31
32 Nydam, & Hare, 2014), the highly consistent patterns of excess suggest that all populations would have
33
34 been affected similarly, enabling interpretations of relative variation within and between populations. The
35
36 admixture suggested by the STRUCTURE analyses also supports the role of introgression in the histo-
37
38 ries of the studied populations. However, mapping of the ddRAD reads to multiple reference genomes
39
40 or to data generated based on different restriction enzymes would be more problematic.
41
42
43

428 **Environmental associations**

44
45
46 Despite the general lack of geographic clustering, there was evidence of local adaptation to the vary-
47
48 ing environments using redundancy analyses, particularly to annual precipitation (Figure 3). Although
49
50 Watson-Jones et al. (2006) found some population structuring within the UK, the same result was not
51
52 found in this study (i.e. no isolation-by-distance within the UK). Furthermore, no evidence of population
53
54 structuring was found in the Spanish populations here, and Maggioni *et al.* (personal communication)
55
56 found no evidence of population structuring in French Atlantic populations. These results could also be
57
58 correlated with annual precipitation; perhaps the strong variation in annual precipitation in the UK (e.g.,
59
60

1
2
3
4 436 a strong west-east gradient) is causing more differentiation between these populations, whereas along
5
6 437 the French range annual precipitation has a smaller gradient. One reason for the importance of annual
7
8 438 precipitation other than water availability could be the influence of precipitation on soil pH. Soil pH is
9
10 439 primarily determined by bedrock, but is also altered by precipitation through leaching of compounds
11
12 440 such as calcium carbonate (Kinzel, 1983). Therefore, although slightly alkaline to neutral soils tend to
13
14 441 form over limestone, secondary acidification can occur under higher precipitation regimes. The soil pH
15
16 442 values recorded here ranged from neutral to strongly acidic (Figure 4). Furthermore, the bedrock of a
17
18 443 large proportion of the populations used here (Table 2) differ from the limestone and chalk cliffs that wild
19
20 444 *B. oleracea* are thought to be predominantly found on (Christensen et al., 2011). For those individuals
21
22 445 where the soil pH was known, the same trend was found here as by Watson-Jones et al. (2006), with
23
24 446 a decrease in plant genetic diversity as soil pH increased (Figure 4). For agriculture and horticulture,
25
26 447 soil pH is an important consideration (Tilman, Balzer, Hill, & Befort, 2011). The change in plant genetic
27
28 448 diversity suggests that soil pH is a strong selective pressure in the wild, causing an adaptive ecolog-
29
30 449 ical bottleneck in locations where it is higher, resulting in lower genetic diversity. These indications of
31
32 450 local adaptation despite a lack of population structure highlight environmental variables that could be
33
34 451 investigated further in wild populations of *B. oleracea*, which regardless of their origin are surviving.

35
36 452 Alongside survival, a huge concern for food security related to climate change is the ability of crop
37
38 453 plants to remain productive under rapidly changing environmental conditions (Lasky et al., 2015). Ob-
39
40 454 taining accurate phenotypic data for adaptive traits is a major barrier as we often do not know the com-
41
42 455 bination of traits that underlie differences in fitness or how these vary with the environment (Kooyers,
43
44 456 Greenlee, Colicchio, Oh, & Blackman, 2015). Although some traits will be locally adaptive due to large
45
46 457 effect loci, the vast majority of adaptive traits are likely to have a polygenic basis (Rockman, 2012),
47
48 458 particularly in the case of multi-trait phenotypes related to environmental gradients. Our results match
49
50 459 these expectations, as no large effect loci were found; however, some were more significantly associ-
51
52 460 ated with the assessed environmental variation than others. The most likely assignment for six of the
53
54 461 top 18 candidate genes was to the receptor-like kinase family (Table 4). This gene family underwent an
55
56 462 expansion that is believed to be a plant-specific adaptation for pathogen defence (Afzal, Wood, & Light-
57
58 463 foot, 2008). Interestingly, Zhang et al. (2014) also found differences in genes related to plant defence
59
60

1
2
3
4 464 when investigating adaptations of rice (*Oryza* sp.) across four continents. These results highlight the
5
6 465 fundamental importance of the immune system to fitness, and suggest that it could be related to envi-
7
8 466 ronmental differences across different spatial scales. Given that immune system genes are among the
9
10 467 best candidates for local adaptation, there is a potential connection between plant genetic diversity, soil
11
12 468 pH and pathogens. It would be interesting to investigate whether less acidic soils host more pathogens,
13
14 469 increasing the selective pressure on the plants and decreasing the plant genetic diversity in these soils.
15
16 470 Overall, the impact of climate change on the spread of virulence of plant pathogens and herbivores, and
17
18 471 the phenological mismatches that may occur between interacting species remain unknown (De Lucia,
19
20 472 Nability, Zavala, & Berenbaum, 2012; Fisher et al., 2012; Yang & Rudolf, 2010). What is clear is that plant
21
22 473 defence will continue to be an important component of crop productivity, warranting further research.

23
24 474 Overall, the results presented here supported the hypothesis that wild populations of *B. oleracea* in
25
26 475 the Atlantic region were established by plants from agricultural and/or horticultural sources. In addition,
27
28 476 regardless of their origin, these wild populations are likely to contain useful genetic resources and should
29
30 477 be considered as valuable populations of a crop wild relative to be investigated further.

31 32 33 34 478 **Acknowledgements**

35
36
37 479 EAM was funded by a University of Glasgow Lord Kelvin Adam Smith PhD studentship; UZI was
38
39 480 funded by a NERC Independent Research Fellowship (NE/L011956); CAC is supported by the BBSRC
40
41 481 (BB/P004202/1); KAM utilised equipment funded by the Wellcome Trust Institutional Strategic Support
42
43 482 Fund (WT097835MF), Wellcome Trust Multi User Equipment Award (WT101650MA) and BBSRC LOLA
44
45 483 award (BB/K003240/1). Part of the work was supported by a British Society for Plant Pathology summer
46
47 484 studentship, and grants from the Botanical Research Fund, and the Blodwen Lloyd Bins trust funded
48
49 485 through the Glasgow Natural History Society. None of the sponsors had any role in the study design,
50
51 486 data collection, analysis, and interpretation or any aspects during the write up and publication of this
52
53 487 work. We thank Danijela Dimitrijević and Deborah Davy for assistance in the field, and Dr Lorenzo
54
55 488 Maggioni for useful discussions.

References

- 489
- 490 Afzal, A. J., Wood, A. J., & Lightfoot, D. A. (2008). Plant receptor-like serine threonine kinases: roles in
491 signaling and plant defense. *Molecular Plant-Microbe Interactions*, *21*(5), 507–517.
- 492 Allender, C., Allainguillaume, J., Lynn, J., & King, G. J. (2007). Simple sequence repeats reveal uneven
493 distribution of genetic diversity in chloroplast genomes of *Brassica oleracea* L. and (n= 9) wild
494 relatives. *Theoretical and Applied Genetics*, *114*(4), 609–618.
- 495 Andrews, K. R., Good, J. M., Miller, M. R., Luikart, G., & Hohenlohe, P. A. (2016). Harnessing the power
496 of RADseq for ecological and evolutionary genomics. *Nature Reviews Genetics*, *17*(2), 81–92.
- 497 BGS. (2018). *British Geological Survey*. <http://www.bgs.ac.uk/>. (Accessed: 2018-08-26)
- 498 BSBI. (2018). *Botanical Society of Britain & Ireland*. <https://bsbi.org>. (Accessed: 2018-08-26)
- 499 Buckley, J., Holub, E. B., Koch, M. A., Vergeer, P., & Mable, B. K. (2018). Restriction associated DNA-
500 genotyping at multiple spatial scales in *Arabidopsis lyrata* reveals signatures of pathogen-mediated
501 selection. *BMC Genomics*, *19*(1), 496.
- 502 Charlesworth, D., Vekemans, X., Castric, V., & Glémin, S. (2005). Plant self-incompatibility systems: a
503 molecular evolutionary perspective. *New Phytologist*, *168*(1), 61–69.
- 504 Christensen, S., von Bothmer, R., Poulsen, G., Maggioni, L., Phillip, M., Andersen, B. A., & Jørgensen,
505 R. B. (2011). AFLP analysis of genetic diversity in leafy kale (*Brassica oleracea* L. convar. *acephala*
506 (DC.) Alef.) landraces, cultivars and wild populations in Europe. *Genetic Resources and Crop*
507 *Evolution*, *58*(5), 657–666.
- 508 Conner, J. K. (2002). Genetic mechanisms of floral trait correlations in a natural population. *Nature*,
509 *420*(6914), 407–410.
- 510 Darwin, C. R. (1859). *The origin of species*. London: John Murray.
- 511 De Lucia, E., Nabity, P., Zavala, J., & Berenbaum, M. (2012). Climate change: resetting plant-insect
512 interactions. *Plant Physiology*, pp–112.
- 513 Diamond, J. (2002). Evolution, consequences and future of plant and animal domestication. *Nature*,
514 *418*(6898), 700.
- 515 Euro+Med PlantBase. (2020). <http://ww2.bgbm.org/EuroPlusMed/>. (Accessed: 2020-05-23)
- 516 Evanno, G., Regnaut, S., & Goudet, J. (2005). Detecting the number of clusters of individuals using the
517 software STRUCTURE: a simulation study. *Molecular Ecology*, *14*(8), 2611–2620.
- 518 Falconer, D. S., & Mackay, T. F. C. (1996). *Introduction to Quantitative Genetics* (4th ed.). Harlow, UK:
519 Longman.
- 520 Fick, S. E., & Hijmans, R. J. (2017). WorldClim 2: new 1-km spatial resolution climate surfaces for global
521 land areas. *International Journal of Climatology*, *37*(12), 4302–4315.
- 522 Fisher, M. C., Henk, D. A., Briggs, C. J., Brownstein, J. S., Madoff, L. C., McCraw, S. L., & Gurr, S. J.
523 (2012). Emerging fungal threats to animal, plant and ecosystem health. *Nature*, *484*(7393), 186.
- 524 Forester, B. R., Lasky, J. R., Wagner, H. H., & Urban, D. L. (2018). Comparing methods for detecting
525 multilocus adaptation with multivariate genotype–environment associations. *Molecular Ecology*,
526 *27*(9), 2215–2233.
- 527 Funk, J. L., Standish, R. J., Stock, W. D., & Valladares, F. (2016). Plant functional traits of dominant
528 native and invasive species in mediterranean-climate ecosystems. *Ecology*, *97*(1), 75–83.
- 529 Gao, Y., Yin, S., Wu, L., Dai, D., Wang, H., Liu, C., & Tang, L. (2017). Genetic diversity and structure of
530 wild and cultivated *Amorphophallus paeoniifolius* populations in southwestern China as revealed
531 by RAD-seq. *Scientific Reports*, *7*(1), 14183.
- 532 Gómez-Campo, C., & Prakash, S. (1999). Origin and domestication. In *Developments in plant genetics*
533 *and breeding* (Vol. 4, pp. 33–58). Elsevier.
- 534 Henriksen, R., Gering, E., & Wright, D. (2018). Feralisation—The understudied counterpoint to domes-
535 tication. In *Origin and evolution of biodiversity* (pp. 183–195). Springer.
- 536 Hijmans, R. J. (2017). geosphere: Spherical Trigonometry [Computer software manual]. Retrieved from
537 <https://CRAN.R-project.org/package=geosphere> (R package version 1.5-7)
- 538 Hoisington, D., Khairallah, M., Reeves, T., Ribaut, J.-M., Skovmand, B., Taba, S., & Warburton, M.
539 (1999). Plant genetic resources: What can they contribute toward increased crop productivity?
540 *Proceedings of the National Academy of Sciences*, *96*(11), 5937–5943.
- 541 Huson, D. H., & Bryant, D. (2005). Application of phylogenetic networks in evolutionary studies. *Molec-
542 ular Biology and Evolution*, *23*(2), 254–267.
- 543 IGME. (2018). *Instituto Geológico y Minero de España*. <http://www.igme.es>. (Accessed: 2018-08-26)

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- 544 llut, D. C., Nydam, M. L., & Hare, M. P. (2014). Defining loci in restriction-based reduced representation
545 genomic data from nonmodel species: sources of bias and diagnostics for optimal clustering.
546 *BioMed Research International*, 2014.
- 547 Kinzel, H. (1983). Influence of limestone, silicates and soil pH on vegetation. In *Physiological plant*
548 *ecology iii* (pp. 201–244). Springer.
- 549 Kitashiba, H., & Nasrallah, J. B. (2014). Self-incompatibility in Brassicaceae crops: lessons for interspe-
550 cific incompatibility. *Breeding Science*, 64(1), 23–37.
- 551 Kofsky, J., Zhang, H., & Song, B.-H. (2018). The untapped genetic reservoir: The past, current, and
552 future applications of the wild soybean (*Glycine soja*). *Frontiers in Plant Science*, 9.
- 553 Kooyers, N. J., Greenlee, A. B., Colicchio, J. M., Oh, M., & Blackman, B. K. (2015). Replicate altitudinal
554 clines reveal that evolutionary flexibility underlies adaptation to drought stress in annual *Mimulus*
555 *guttatus*. *New Phytologist*, 206(1), 152–165.
- 556 Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods*,
557 9(4), 357.
- 558 Lanner-Herrera, C., Gustafeson, M., Filt, A., & Bryngelsson, T. (1996). Diversity in natural populations
559 of wild *Brassica oleracea* as estimated by isozyme and RAPD analysis. *Genetic Resources and*
560 *Crop Evolution*, 43(1), 13–23.
- 561 Lasky, J. R., Upadhyaya, H. D., Ramu, P., Deshpande, S., Hash, C. T., Bonnette, J., ... others (2015).
562 Genome-environment associations in sorghum landraces predict adaptive traits. *Science Ad-*
563 *vances*, 1(6), e1400218.
- 564 Lázaro, A., & Aguinagalde, I. (1998). Genetic diversity in *Brassica oleracea* L.(Cruciferae) and wild
565 relatives (2 n= 18) using isozymes. *Annals of Botany*, 82(6), 821–828.
- 566 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., ... Durbin, R. (2009). The sequence
567 alignment/map format and SAMtools. *Bioinformatics*, 25(16), 2078–2079.
- 568 Liu, S., Liu, Y., Yang, X., Tong, C., Edwards, D., Parkin, I. A., ... others (2014). The *Brassica oleracea*
569 genome reveals the asymmetrical evolution of polyploid genomes. *Nature Communications*, 5.
- 570 Maggioni, L. (2015). *Domestication of Brassica oleracea l.* (Unpublished doctoral dissertation). Swedish
571 University of Agricultural Studies, Alnarp, Sweden.
- 572 Maggioni, L., von Bothmer, R., Poulsen, G., & Lipman, E. (2018). Domestication, diversity and use
573 of *Brassica oleracea* L., based on ancient Greek and Latin texts. *Genetic Resources and Crop*
574 *Evolution*, 65(1), 137–159.
- 575 Marandel, F., Charrier, G., Lamy, J.-B., Le Cam, S., Lorange, P., & Trenkel, V. M. (2020). Estimating
576 effective population size using RADseq: Effects of SNP selection and sample size. *Ecology and*
577 *Evolution*, 10(4), 1929–1937.
- 578 Martin, M. (2011, May). Cutadapt removes adapter sequences from high-throughput sequencing
579 reads. *EMBnet.journal*, 17(1), 10–12. Retrieved from [http://journal.embnet.org/index.php/](http://journal.embnet.org/index.php/embnetjournal/article/view/200)
580 [embnetjournal/article/view/200](http://journal.embnet.org/index.php/embnetjournal/article/view/200) doi: <http://dx.doi.org/10.14806/ej.17.1.200>
- 581 Milla, R., Osborne, C. P., Turcotte, M. M., & Violle, C. (2015). Plant domestication through an ecological
582 lens. *Trends in Ecology & Evolution*.
- 583 Mitchell, N. (1976). The Status of *Brassica oleracea* L. subsp. *oleracea* (wild cabbage) in the British
584 Isles. *Watsonia*, 11, 97–103.
- 585 Nazareno, A. G., Bemmels, J. B., Dick, C. W., & Lohmann, L. G. (2017). Minimum sample sizes for
586 population genomics: An empirical study from an Amazonian plant species. *Molecular Ecology*
587 *Resources*, 17(6), 1136–1147.
- 588 Nei, M. (1973). Analysis of gene diversity in subdivided populations. *Proceedings of the National*
589 *Academy of Sciences*, 70(12), 3321–3323.
- 590 Oksanen, J., Blanchet, F. G., Friendly, M., Kindt, R., Legendre, P., McGlenn, D., ... Wagner, H. (2017).
591 vegan: Community Ecology Package [Computer software manual]. Retrieved from [https://CRAN](https://CRAN.R-project.org/package=vegan)
592 [.R-project.org/package=vegan](https://CRAN.R-project.org/package=vegan) (R package version 2.4-3)
- 593 Paris, J. R., Stevens, J. R., & Catchen, J. M. (2017). Lost in parameter space: a road map for stacks.
594 *Methods in Ecology and Evolution*, 8(10), 1360–1373.
- 595 Prentis, P. J., Wilson, J. R., Dormontt, E. E., Richardson, D. M., & Lowe, A. J. (2008). Adaptive evolution
596 in invasive species. *Trends in Plant Science*, 13(6), 288–294.
- 597 Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of population structure using multilocus
598 genotype data. *Genetics*, 155(2), 945–959.

- 1
2
3
4
599 Purugganan, M. D., & Fuller, D. Q. (2009). The nature of selection during plant domestication. *Nature*,
600 457(7231), 843–848.
- 601 Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features.
602 *Bioinformatics*, 26(6), 841–842.
- 603 R Core Team. (2017). R: A language and environment for statistical computing [Computer software
604 manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- 605 Rauf, S., Teixeira da Silva, J., Khan, A. A., & Naveed, A. (2010). Consequences of plant breeding on
606 genetic diversity. *International Journal of Plant Breeding*, 4(1), 1–21.
- 607 Raybould, A., Mogg, R., Clarke, R., Gliddon, C., & Gray, A. (1999). Variation and population structure
608 at microsatellite and isozyme loci in wild cabbage (*Brassica oleracea* L.) in Dorset (UK). *Genetic
609 Resources and Crop Evolution*, 46(4), 351–360.
- 610 Revelle, W. (2018). psych: Procedures for Psychological, Psychometric, and Personality Research
611 [Computer software manual]. Evanston, Illinois. Retrieved from [https://CRAN.R-project.org/
612 package=psych](https://CRAN.R-project.org/package=psych) (R package version 1.8.10)
- 613 Rochette, N. C., & Catchen, J. M. (2017). Deriving genotypes from RAD-seq short-read data using
614 Stacks. *Nature Protocols*, 12(12), 2640.
- 615 Rockman, M. V. (2012). The QTN program and the alleles that matter for evolution: all that's gold does
616 not glitter. *Evolution*, 66(1), 1–17.
- 617 Rousset, F., & Raymond, M. (1995). Testing heterozygote excess and deficiency. *Genetics*, 140(4),
618 1413–1419.
- 619 Sánchez-Yélamo, M. (2014). Characterisation of wild cabbage (*Brassica oleracea* L.) based on iso-
620 zyme data. considerations on the current status of this taxon in Spain. *Genetic Resources and
621 Crop Evolution*, 61(7), 1295–1306.
- 622 Scheepens, J., Frei, E. S., Armbruster, G. F., & Stöcklin, J. (2012). Pollen dispersal and gene flow within
623 and into a population of the alpine monocarpic plant *Campanula thyrsoidea*. *Annals of Botany*,
624 110(7), 1479–1488.
- 625 Schranz, M. E., Lysak, M. A., & Mitchell-Olds, T. (2006). The ABC's of comparative genomics in the
626 Brassicaceae: building blocks of crucifer genomes. *Trends in Plant Science*, 11(11), 535–542.
- 627 Sharbel, T. F., Haubold, B., & Mitchell-Olds, T. (2000). Genetic isolation by distance in *Arabidopsis
628 thaliana*: biogeography and postglacial colonization of Europe. *Molecular Ecology*, 9(12), 2109–
629 2118.
- 630 Slatkin, M. (1987). Gene flow and the geographic structure of natural populations. *Science*, 236(4803),
631 787–792.
- 632 Song, K., Osborn, T. C., & Williams, P. H. (1990). *Brassica* taxonomy based on nuclear restriction
633 fragment length polymorphisms (RFLPs). *Theoretical and Applied Genetics*, 79(4), 497–506.
- 634 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
635 phylogenies. *Bioinformatics*, 30(9), 1312–1313.
- 636 Stansell, Z., Hyma, K., Fresnedo-Ramírez, J., Sun, Q., Mitchell, S., Björkman, T., & Hua, J. (2018).
637 Genotyping-by-sequencing of *Brassica oleracea* vegetables reveals unique phylogenetic patterns,
638 population structure and domestication footprints. *Horticulture Research*, 5.
- 639 Tilman, D., Balzer, C., Hill, J., & Befort, B. L. (2011). Global food demand and the sustainable intensifi-
640 cation of agriculture. *Proceedings of the National Academy of Sciences*, 108(50), 20260–20264.
- 641 USDA, D. o. A. (1998). *Soil Quality Indicators: pH*. [https://www.nrcs.usda.gov/Internet/
642 FSE_DOCUMENTS/nrcs142p2_052208.pdf](https://www.nrcs.usda.gov/Internet/FSE_DOCUMENTS/nrcs142p2_052208.pdf). (Accessed: 2016-04-20)
- 643 von Wettberg, E. J., Chang, P. L., Başdemir, F., Carrasquilla-Garcia, N., Korbu, L. B., Moenga, S. M., ...
644 others (2018). Ecology and genomics of an important crop wild relative as a prelude to agricultural
645 innovation. *Nature Communications*, 9(1), 649.
- 646 Walley, P. G., Teakle, G. R., Moore, J. D., Allender, C. J., Pink, D. A., Buchanan-Wollaston, V., & Barker,
647 G. C. (2012). Developing genetic resources for pre-breeding in *Brassica oleracea* L.: an overview
648 of the UK perspective. *Journal of Plant Biotechnology*, 39(1), 62–68.
- 649 Watson-Jones, S., Maxted, N., & Ford-Lloyd, B. (2006). Population baseline data for monitoring genetic
650 diversity loss for 2010: a case study for *Brassica* species in the UK. *Biological Conservation*,
651 132(4), 490–499.
- 652 Wichmann, M. C., Alexander, M. J., Soons, M. B., Galsworthy, S., Dunne, L., Gould, R., ... Bullock,
653 J. M. (2009). Human-mediated dispersal of seeds over long distances. *Proceedings of the Royal
654 Society of London B: Biological Sciences*, 276(1656), 523–532.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- 655 Wu, Z., Wang, B., Chen, X., Wu, J., King, G. J., Xiao, Y., & Liu, K. (2016). Evaluation of linkage
656 disequilibrium pattern and association study on seed oil content in *Brassica napus* using ddRAD
657 sequencing. *PLoS One*, *11*(1), e0146383.
- 658 Xiao, Z., Han, F., Hu, Y., Xue, Y., Fang, Z., Yang, L., ... others (2019). Overcoming cabbage crossing
659 incompatibility by the development and application of self-compatibility-QTL-specific markers and
660 genome-wide background analysis. *Frontiers in Plant Science*, *10*, 189.
- 661 Yang, L. H., & Rudolf, V. (2010). Phenology, ontogeny and the effects of climate change on the timing
662 of species interactions. *Ecology Letters*, *13*(1), 1–10.
- 663 Yousef, E. A., Mueller, T., Börner, A., & Schmid, K. J. (2018). Comparative analysis of genetic diver-
664 sity and differentiation of cauliflower (*Brassica oleracea* var. botrytis) accessions from two *ex situ*
665 genebanks. *PLoS One*, *13*(2), e0192062.
- 666 Yu, J., Zhao, M., Wang, X., Tong, C., Huang, S., Tehrim, S., ... Liu, S. (2013). Bolbase: a comprehensive
667 genomics database for *Brassica oleracea*. *BMC Genomics*, *14*(1), 664.
- 668 Zhang, Q.-J., Zhu, T., Xia, E.-H., Shi, C., Liu, Y.-L., Zhang, Y., ... others (2014). Rapid diversification of
669 five *Oryza* AA genomes associated with rice adaptation. *Proceedings of the National Academy of*
670 *Sciences*, *111*(46), E4954–E4962.

671 **Data Accessibility**

672 **The sequencing data and associated meta data are available on the European Nucleotide Archive under**
673 **the study accession number: PRJEB38464 (<http://www.ebi.ac.uk/ena/data/view/PRJEB38464>)**

674 **Author Contributions**

675 EAM and BKM conceived and proposed the study to the co-authors; EAM, BKM, CAC, and UZI then
676 developed the study design. BKM, CAC, and UZI obtained the main funding award. EAM and EAK
677 collected the data. EAM carried out statistical analysis, and wrote the paper. KAM carried out the
678 double-digest RAD-sequencing. All authors reviewed the final version of the manuscript and agreed to
679 its content before submission.

Table 1: Estimates of genetic diversity within wild *B. oleracea* populations from previous studies using different molecular markers. H_E is expected heterozygosity estimated using Nei's gene diversity (Nei, 1973).

Study	Molecular marker	H_E	Populations
Lanner-Herrera et al. (1996)	Isozymes	0.10 – 0.56	France, Spain, UK
Lázaro and Aguinagalde (1998)	Isozymes	0.26 – 0.30	France, Spain, UK
Raybould et al. (1999)	Isozymes	0.40 (0.18 – 0.41)†	UK
Raybould et al. (1999)	Microsatellites	0.36 (0.21 – 0.33)†	UK
Watson-Jones et al. (2006)	AFLPs	0.19 – 0.33	UK
Christensen et al. (2011)	AFLPs	0.23, 0.20	Spain, UK
Maggioni <i>et al.</i> (pers. comm. 2019)	AFLPs	0.25	France

† – pooled population H_E with the range of estimates from individual populations shown in brackets.

Table 2: A summary of the natural populations of *B. oleracea* used in this study, including: the bedrock, the first time the population was recorded, the number of individuals sequenced, and the number of individuals included in down-stream analyses.

Region	Population	Bedrock†	First population record‡	Number sequenced	Number included§
ES	Auchmithie	Red basic sandstone	1913	4	4
ES	Crail	Sandstone & limestone	1840	4	4
ES	Fortrose	Sandstone	1968	4	3
WS	Kildonan Castle	Sandstone & limestone	1987	4	4
NEE	Tynemouth	Sandstone & limestone	1805	4	4
NEE	Staithes	Shale & sandstone	1831	4	–
NEE	Whitby	Shale	1906	4	4
NW	Little Orme	Limestone	1895	4	–
SW	Tenby	Siltstone & sandstone	1773	4	4
SW	Llantwit Major	Limestone	1850	4	4
SWE	Prussia Cove	Slate, shale & siltstone	1871	4	4
SWE	Fowey	Shale & siltstone	1805	4	4
SWE	West Looe	Siltstone & sandstone	1971	4	2
SWE	St. Aldhelm's Head	Limestone	1933	4	4
A	Cudillero	Slate & sandstone		4	4
A	Playa de Xágo	Sandstone & Dolomite		4	4
A	Cabo de Peñas	Slate & quartzite		4	4
A	Playas de Viodo	Slate & shale		4	4
A	Tazonas	Dolomite & limestone		4	–
C	Playa Pedrero	Quartzites		4	4
C	La Franca	Quartzites		4	–
BC	San Juan de Gaxtelugatxe	Limestone		4	4
BC	Getaría	Limestone		4	4
BC	San Sebastian	Calcareous sandstone		4	3
Total:				96	76

† data obtained from the British Geological Survey (<https://www.bgs.ac.uk>) and the Instituto Geológico y Minero de España (<http://www.igme.es>). Region codes: ES – East Scotland, WS – West Scotland, NEE – North-eastern England, NW – North Wales, SW – South Wales, SWE – South-western England, A – Asturias Spain, C – Cantabrica Spain, BC – Basque Country Spain. ‡ data obtained from the Botanical Society of Britain & Ireland (<https://bsbi.org>). § indicates where data was lost in quality filtering of sequences and not included in down-stream analyses.

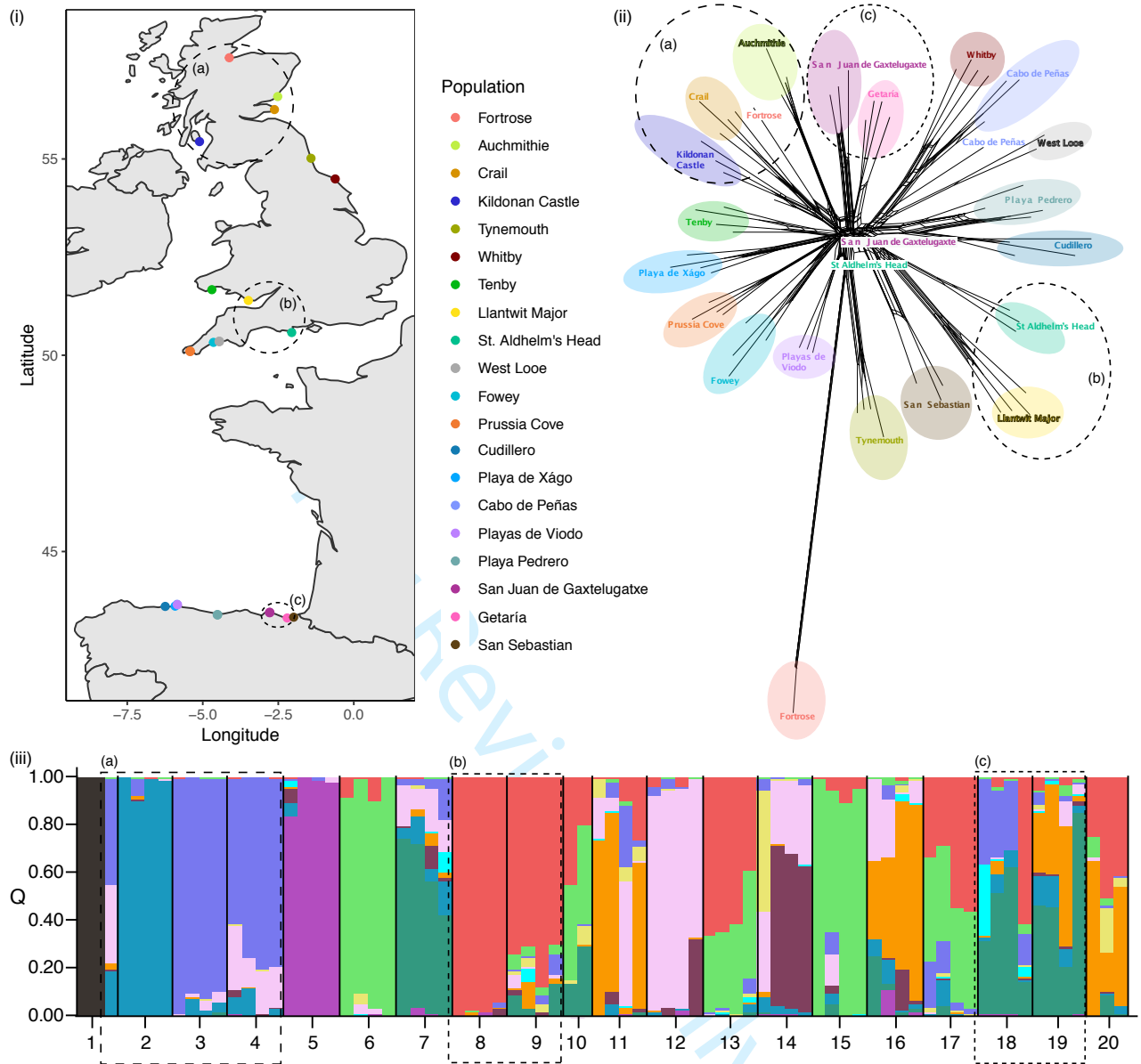


Figure 1: Population structuring of wild populations of *B. oleracea*. (i) Location of the populations considered here. (ii) Clustering of samples from RAXML (v8.2; GTRCAT model and 1000 maximum likelihood bootstrap replicates), visualised in SplitsTree4. (iii) STRUCTURE plot illustrating shared genetic ancestry for $K = 12$, ordered by population: 1 – Fortrose; 2 – Auchmithie; 3 – Craíl; 4 – Kildonan Castle; 5 – Tynemouth; 6 – Whitby; 7 – Tenby; 8 – Llantwit Major; 9 – St. Aldhelm's Head; 10 – West Looe; 11 – Fowey; 12 – Prussia Cova; 13 – Cudillero; 14 – Playa de Xágo; 15 – Cabo de Peñas; 16 – Playas de Viodo; 17 – Playa Pedrero; 18 – San Juan de Gaxtelugatxe; 19 – Getaría; 20 – San Sebastian. Across the figures the same colours and numbering is used for each population. The dashed lines and letters indicate some clustering: (a) populations in Scotland; (b) populations closest to the Welsh-English border; and (c) populations in the Basque Country, Spain (excluding San Sebastian).

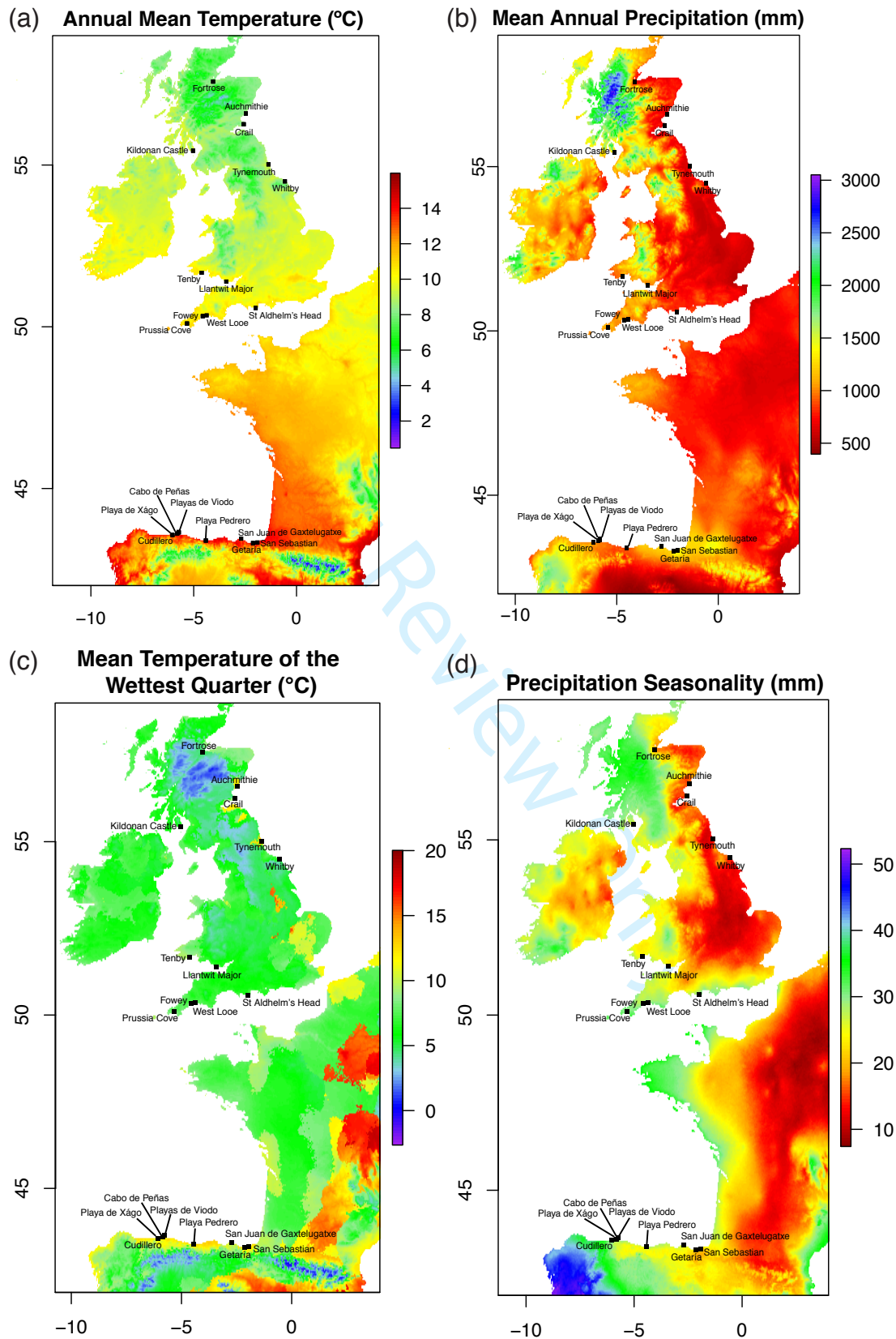
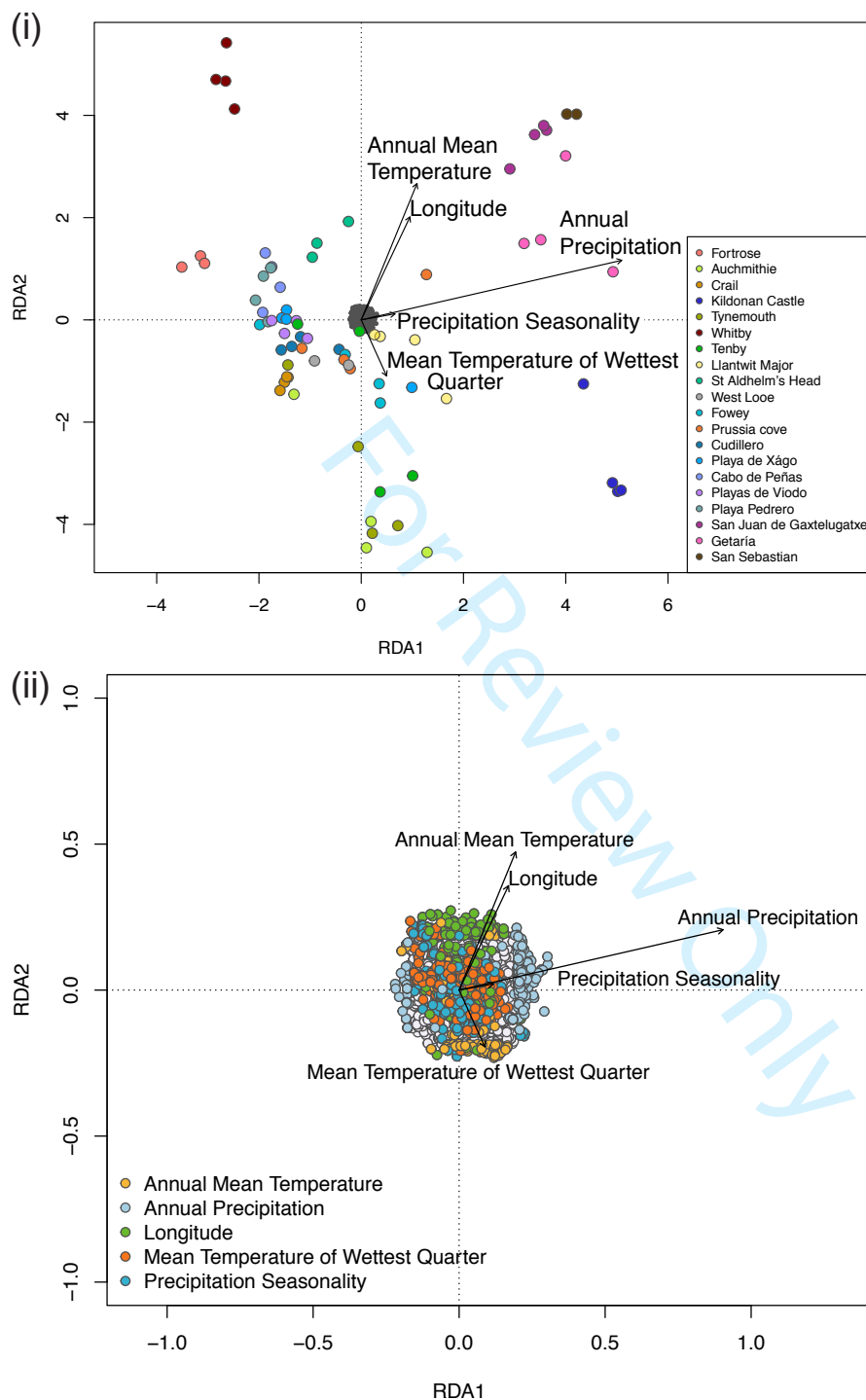


Figure 2: The distribution of sampled populations in relation to various climate variables: (a) annual mean temperature (°C); (b) mean annual precipitation (mm); (c) mean temperature of wettest quarter (°C); (d) precipitation seasonality (mm). These are averages between 1970 – 2000 obtained from the WorldClim database (Fick & Hijmans, 2017).

Table 3: Summary statistics of within *B. oleracea* population genetic diversity based on both variant nucleotide sites alone (var) and all sites (all) from dataset 1, showing: the number of individuals (N), the number of private alleles (PRI), expected heterozygosity (H_E), observed heterozygosity (H_O) and percentage of polymorphic loci (%).

Region†	Population	N	PRI	H_E		H_O		% All
				Var	All	Var	All	
ES	Authmithie	4	1683	0.1043	0.0012	0.1202	0.0014	33.0
ES	Crail	4	1727	0.1327	0.0019	0.1267	0.0018	52.8
ES	Fortrose	3	12951	0.2006	0.0032	0.1962	0.0031	76.4
WS	Kildonan Castle	4	1014	0.0903	0.0014	0.0944	0.0014	40.8
NEE	Tynemouth	4	1476	0.1023	0.0013	0.0881	0.0011	36.4
NEE	Whitby	4	1573	0.1200	0.0020	0.1184	0.0020	56.7
SW	Tenby	4	1568	0.1227	0.0014	0.1153	0.0013	40.5
SW	Llantwit Major	4	2073	0.1390	0.0023	0.1231	0.0022	66.2
SWE	Prussia Cove	4	1454	0.1019	0.0016	0.1064	0.0017	45.5
SWE	Fowey	4	1137	0.1126	0.0018	0.1083	0.0017	53.4
SWE	West Looe	2	1412	0.1150	0.0011	0.1328	0.0013	27.1
SWE	St. Aldhelm's Head	4	2470	0.1486	0.0014	0.1676	0.0016	39.4
A	Cudillero	4	716	0.0918	0.0015	0.0938	0.0016	44.3
A	Playa de Xágo	4	1583	0.1140	0.0012	0.1191	0.0012	33.4
A	Cabo de Peñas	4	698	0.0933	0.0015	0.0910	0.0014	42.5
A	Playas de Viodo	4	503	0.0545	0.0004	0.0580	0.0004	11.2
C	Playa Pedrero	4	1741	0.1313	0.0014	0.1408	0.0015	38.5
BC	San Juan de Gaxtelugatxe	4	2608	0.1423	0.0012	0.1471	0.0012	34.0
BC	Getaría	4	1550	0.1280	0.0021	0.1391	0.0023	59.8
BC	San Sebastian	3	2516	0.1530	0.0023	0.1538	0.0023	61.4

†Region codes: ES – East Scotland, WS – West Scotland, NEE – North-eastern England, SW – South Wales, SWE – South-western England, A – Asturias Spain, C – Cantabrica Spain, BC – Basque Country Spain.



53
54
55
56
57
58
59
60

Figure 3: (i) Redundancy analysis (RDA) ordination plot of the association between *B. oleracea* individuals (coloured points) and SNPs (dark grey points), with environmental variables. The different colours indicate which population each individual was from. (ii) RDA ordination plot of the SNPs alone, coloured for the environmental variable with which they were most strongly associated. For both (i) & (ii) the arrows indicate the environmental predictors and the strength of the association.

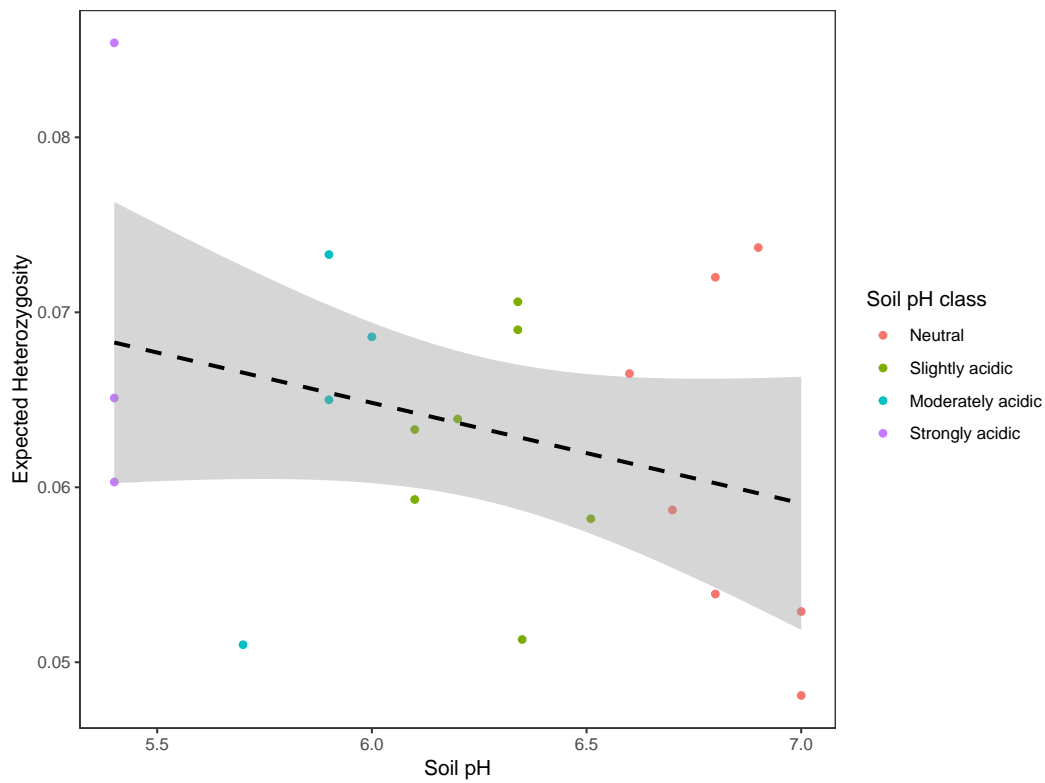


Figure 4: The relationship between expected heterozygosity and soil pH for 21 individuals from four soil pH classes categorised into: Neutral (6.6 - 7.3), Slightly acidic (6.1 - 6.5), Moderately acidic (5.6- 6.0) and Strongly acidic (5.0 - 5.5) based on USDA (1998). A linear model was used to fit a regression line (dashed black line), the standard error is shown in grey, p-value > 0.05.

Table 4: The top 18 candidate SNPs that mapped to unique genes in the *B. oleracea* reference genome and their annotations from 'Bolbase' (Yu et al., 2013).

Chromosome	Location	Identity	X	Bolbase gene name	Potential protein	Function
C09	32879582	1	-	Bol019890	Ribonucleotide reductase-related	Fatty acid metabolic process, creation of DNA from RNA
C04	39737611	0.999979	-	Bol021601	Unknown	
C09	8499546	1	+	Bol032146	Basic helix-loop-helix dimerisation region	Nucleus transcription regulation
C07	43014116	1	-	Bol042101	Toll-Interleukin receptor	Signal transduction, immune response, disease resistance
C02	233586	1	+	Bol012817	Laccase/multicopper oxidase	Copper ion binding, metabolic process, maybe formation and degradation of lignin
C04	22051514	0.999656	+	Bol044300	Protein kinase - serine/threonine	Protein kinase activity, signalling, plant defence
C03	29308196	0.472347	-	Bol012462	PIK-related kinase	Binding and DNA repair
C03	48963472	0.99438	+	Bol029900	Protein kinase	Protein kinase activity, signalling, plant defence
C04	28456859	0.999661	-	Bol009961	Cystathionine beta-synthase	Vitamin B6 pathway?
C03	9456274	1	-	Bol005573	Unknown	
C05	2317477	0.580051	-	Bol041075	Pentotricopeptide repeat	Often essential in mitochondria
C04	35972614	0.304057	+	Bol037830	Bacterial transferase haxapeptide repeat	Binding and transferase activity
C04	35104965	0.996501	+	Bol037950	Cyclin-like F-box	Growth and development
C03	2461137	0.999261	-	Bol034275	Serine/threonine-protein kinase	Signalling, plant defence
C02	233586	0.168963	-	Bol012816	Serine/threonine-protein kinase	Signalling, plant defence
C01	11164295	0.999978	+	Bol039465	Initiation factor eIF-4 gamma, MA3	
C01	11431159	1	+	Bol039505	Heat shock protein Hsp20	
C01	12106862	0.918256	-	Bol039585	F-box associated	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



1143x1524mm (72 x 72 DPI)

Mittell_et_al 2020 Supplementary Information output from populations in stacks for the three datasets that were analysed in the manuscript.

Dataset 1 -- Within individuals. All individuals analysed together as one population.

Dataset 2 -- As dataset 1, but only one SNP was retained randomly per RAD-locus.

Dataset 3 -- Within populations. Individuals were assigned to the population they were sampled from.

Populations log for dataset 1

Distribution of population loci.

Distribution of valid loci matched to catalog locus.

Valid samples at locus Count

1 313602

2 65350

3 18307

4 8690

5 5498

6 4104

7 3142

8 2794

9 2283

10 2055

11 1949

12 1761

13 1652

14 1547

15 1455

16 1439

17 1342

18 1296

19 1283

20 1190

21 1229

22 1175

23 1156

24 1093

25 1088

26 1094

27 1030

28 1003

29 1016

30 960

31 967

32 918

33 950

34 990

For Review Only

1		
2		
3	35	977
4	36	933
5	37	830
6	38	859
7	39	910
8	40	898
9	41	838
10	42	848
11	43	830
12	44	852
13	45	869
14	46	868
15	47	804
16	48	817
17	49	888
18	50	827
19	51	849
20	52	880
21	53	778
22	54	812
23	55	891
24	56	838
25	57	873
26	58	840
27	59	831
28	60	906
29	61	947
30	62	958
31	63	952
32	64	913
33	65	1069
34	66	1128
35	67	1178
36	68	1243
37	69	1275
38	70	1335
39	71	1237
40	72	1163
41	73	998
42	74	626
43	75	400
44	76	326

Distribution of confounded loci at catalog locus.

Confounded samples at locus Count

0 493502

Distribution of missing loci at catalog loci.

56
57
58
59
60

1
2
3 # Absent samples at locus Count
4 0 493502
5 # Distribution of population loci after applying locus constraints.
6 # Distribution of valid loci matched to catalog locus.
7 # Valid samples at locus Count
8 45 671
9 46 710
10 47 735
11 48 703
12 49 681
13 50 721
14 51 713
15 52 640
16 53 677
17 54 681
18 55 737
19 56 711
20 57 722
21 58 721
22 59 768
23 60 858
24 61 750
25 62 765
26 63 691
27 64 699
28 65 667
29 66 582
30 67 516
31 68 433
32 69 305
33 70 241
34 71 171
35 72 140
36 73 120
37 74 115
38 75 151
39 76 145
40 # Distribution of confounded loci at catalog locus.
41 # Confounded samples at locus Count
42 0 17940
43 # Distribution of missing loci at catalog loci.
44 # Absent samples at locus Count
45 0 309
46 1 282
47 2 352
48 3 547
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only

1		
2		
3	4	847
4	5	1123
5	6	1390
6	7	1588
7	8	1603
8	9	1763
9	10	1710
10	11	1501
11	12	1278
12	13	1122
13	14	858
14	15	598
15	16	425
16	17	289
17	18	173
18	19	96
19	20	51
20	21	21
21	22	9
22	23	4
23	25	1

27 Population 1 contained 0 incompatible loci -- more than two alleles present.

28 Population 2 contained 0 incompatible loci -- more than two alleles present.

29 Population 3 contained 0 incompatible loci -- more than two alleles present.

30 Population 4 contained 0 incompatible loci -- more than two alleles present.

31 Population 5 contained 0 incompatible loci -- more than two alleles present.

32 Population 6 contained 0 incompatible loci -- more than two alleles present.

33 Population 7 contained 0 incompatible loci -- more than two alleles present.

34 Population 8 contained 0 incompatible loci -- more than two alleles present.

35 Population 9 contained 0 incompatible loci -- more than two alleles present.

36 Population 10 contained 0 incompatible loci -- more than two alleles present.

37 Population 11 contained 0 incompatible loci -- more than two alleles present.

38 Population 12 contained 0 incompatible loci -- more than two alleles present.

39 Population 13 contained 0 incompatible loci -- more than two alleles present.

40 Population 14 contained 0 incompatible loci -- more than two alleles present.

41 Population 15 contained 0 incompatible loci -- more than two alleles present.

42 Population 16 contained 0 incompatible loci -- more than two alleles present.

43 Population 17 contained 0 incompatible loci -- more than two alleles present.

44 Population 18 contained 0 incompatible loci -- more than two alleles present.

45 Population 19 contained 0 incompatible loci -- more than two alleles present.

46 Population 20 contained 0 incompatible loci -- more than two alleles present.

47 Population 21 contained 0 incompatible loci -- more than two alleles present.

48 Population 22 contained 0 incompatible loci -- more than two alleles present.

49 Population 25 contained 0 incompatible loci -- more than two alleles present.

50 Population 26 contained 0 incompatible loci -- more than two alleles present.

51 Population 30 contained 0 incompatible loci -- more than two alleles present.

52

53

54

55

56

57

58

59

60

1
2
3 Population 31 contained 0 incompatible loci -- more than two alleles present.
4 Population 32 contained 0 incompatible loci -- more than two alleles present.
5 Population 33 contained 0 incompatible loci -- more than two alleles present.
6 Population 34 contained 0 incompatible loci -- more than two alleles present.
7 Population 35 contained 0 incompatible loci -- more than two alleles present.
8 Population 36 contained 0 incompatible loci -- more than two alleles present.
9 Population 37 contained 0 incompatible loci -- more than two alleles present.
10 Population 38 contained 0 incompatible loci -- more than two alleles present.
11 Population 39 contained 0 incompatible loci -- more than two alleles present.
12 Population 40 contained 0 incompatible loci -- more than two alleles present.
13 Population 41 contained 0 incompatible loci -- more than two alleles present.
14 Population 45 contained 0 incompatible loci -- more than two alleles present.
15 Population 46 contained 0 incompatible loci -- more than two alleles present.
16 Population 47 contained 0 incompatible loci -- more than two alleles present.
17 Population 48 contained 0 incompatible loci -- more than two alleles present.
18 Population 49 contained 0 incompatible loci -- more than two alleles present.
19 Population 50 contained 0 incompatible loci -- more than two alleles present.
20 Population 51 contained 0 incompatible loci -- more than two alleles present.
21 Population 52 contained 0 incompatible loci -- more than two alleles present.
22 Population 53 contained 0 incompatible loci -- more than two alleles present.
23 Population 54 contained 0 incompatible loci -- more than two alleles present.
24 Population 55 contained 0 incompatible loci -- more than two alleles present.
25 Population 56 contained 0 incompatible loci -- more than two alleles present.
26 Population 57 contained 0 incompatible loci -- more than two alleles present.
27 Population 58 contained 0 incompatible loci -- more than two alleles present.
28 Population 59 contained 0 incompatible loci -- more than two alleles present.
29 Population 60 contained 0 incompatible loci -- more than two alleles present.
30 Population 61 contained 0 incompatible loci -- more than two alleles present.
31 Population 62 contained 0 incompatible loci -- more than two alleles present.
32 Population 63 contained 0 incompatible loci -- more than two alleles present.
33 Population 64 contained 0 incompatible loci -- more than two alleles present.
34 Population 65 contained 0 incompatible loci -- more than two alleles present.
35 Population 66 contained 0 incompatible loci -- more than two alleles present.
36 Population 68 contained 0 incompatible loci -- more than two alleles present.
37 Population 69 contained 0 incompatible loci -- more than two alleles present.
38 Population 70 contained 0 incompatible loci -- more than two alleles present.
39 Population 71 contained 0 incompatible loci -- more than two alleles present.
40 Population 72 contained 0 incompatible loci -- more than two alleles present.
41 Population 73 contained 0 incompatible loci -- more than two alleles present.
42 Population 74 contained 0 incompatible loci -- more than two alleles present.
43 Population 75 contained 0 incompatible loci -- more than two alleles present.
44 Population 76 contained 0 incompatible loci -- more than two alleles present.
45 Population 77 contained 0 incompatible loci -- more than two alleles present.
46 Population 86 contained 0 incompatible loci -- more than two alleles present.
47 Population 87 contained 0 incompatible loci -- more than two alleles present.
48 Population 88 contained 0 incompatible loci -- more than two alleles present.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Population 89 contained 0 incompatible loci -- more than two alleles present.
4 Population 90 contained 0 incompatible loci -- more than two alleles present.
5 Population 91 contained 0 incompatible loci -- more than two alleles present.
6 Population 92 contained 0 incompatible loci -- more than two alleles present.
7 Population 93 contained 0 incompatible loci -- more than two alleles present.
8 Population 1 contained 0 incompatible loci -- more than two alleles present.
9 Population 2 contained 0 incompatible loci -- more than two alleles present.
10 Population 3 contained 0 incompatible loci -- more than two alleles present.
11 Population 4 contained 0 incompatible loci -- more than two alleles present.
12 Population 5 contained 0 incompatible loci -- more than two alleles present.
13 Population 6 contained 0 incompatible loci -- more than two alleles present.
14 Population 7 contained 0 incompatible loci -- more than two alleles present.
15 Population 8 contained 0 incompatible loci -- more than two alleles present.
16 Population 9 contained 0 incompatible loci -- more than two alleles present.
17 Population 10 contained 0 incompatible loci -- more than two alleles present.
18 Population 11 contained 0 incompatible loci -- more than two alleles present.
19 Population 12 contained 0 incompatible loci -- more than two alleles present.
20 Population 13 contained 0 incompatible loci -- more than two alleles present.
21 Population 14 contained 0 incompatible loci -- more than two alleles present.
22 Population 15 contained 0 incompatible loci -- more than two alleles present.
23 Population 16 contained 0 incompatible loci -- more than two alleles present.
24 Population 17 contained 0 incompatible loci -- more than two alleles present.
25 Population 18 contained 0 incompatible loci -- more than two alleles present.
26 Population 19 contained 0 incompatible loci -- more than two alleles present.
27 Population 20 contained 0 incompatible loci -- more than two alleles present.
28 Population 21 contained 0 incompatible loci -- more than two alleles present.
29 Population 22 contained 0 incompatible loci -- more than two alleles present.
30 Population 25 contained 0 incompatible loci -- more than two alleles present.
31 Population 26 contained 0 incompatible loci -- more than two alleles present.
32 Population 30 contained 0 incompatible loci -- more than two alleles present.
33 Population 31 contained 0 incompatible loci -- more than two alleles present.
34 Population 32 contained 0 incompatible loci -- more than two alleles present.
35 Population 33 contained 0 incompatible loci -- more than two alleles present.
36 Population 34 contained 0 incompatible loci -- more than two alleles present.
37 Population 35 contained 0 incompatible loci -- more than two alleles present.
38 Population 36 contained 0 incompatible loci -- more than two alleles present.
39 Population 37 contained 0 incompatible loci -- more than two alleles present.
40 Population 38 contained 0 incompatible loci -- more than two alleles present.
41 Population 39 contained 0 incompatible loci -- more than two alleles present.
42 Population 40 contained 0 incompatible loci -- more than two alleles present.
43 Population 41 contained 0 incompatible loci -- more than two alleles present.
44 Population 45 contained 0 incompatible loci -- more than two alleles present.
45 Population 46 contained 0 incompatible loci -- more than two alleles present.
46 Population 47 contained 0 incompatible loci -- more than two alleles present.
47 Population 48 contained 0 incompatible loci -- more than two alleles present.
48 Population 49 contained 0 incompatible loci -- more than two alleles present.
49
50
51
52
53
54
55
56
57
58
59
60

1
 2
 3 Population 50 contained 0 incompatible loci -- more than two alleles present.
 4 Population 51 contained 0 incompatible loci -- more than two alleles present.
 5 Population 52 contained 0 incompatible loci -- more than two alleles present.
 6 Population 53 contained 0 incompatible loci -- more than two alleles present.
 7 Population 54 contained 0 incompatible loci -- more than two alleles present.
 8 Population 55 contained 0 incompatible loci -- more than two alleles present.
 9 Population 56 contained 0 incompatible loci -- more than two alleles present.
 10 Population 57 contained 0 incompatible loci -- more than two alleles present.
 11 Population 58 contained 0 incompatible loci -- more than two alleles present.
 12 Population 59 contained 0 incompatible loci -- more than two alleles present.
 13 Population 60 contained 0 incompatible loci -- more than two alleles present.
 14 Population 61 contained 0 incompatible loci -- more than two alleles present.
 15 Population 62 contained 0 incompatible loci -- more than two alleles present.
 16 Population 63 contained 0 incompatible loci -- more than two alleles present.
 17 Population 64 contained 0 incompatible loci -- more than two alleles present.
 18 Population 65 contained 0 incompatible loci -- more than two alleles present.
 19 Population 66 contained 0 incompatible loci -- more than two alleles present.
 20 Population 68 contained 0 incompatible loci -- more than two alleles present.
 21 Population 69 contained 0 incompatible loci -- more than two alleles present.
 22 Population 70 contained 0 incompatible loci -- more than two alleles present.
 23 Population 71 contained 0 incompatible loci -- more than two alleles present.
 24 Population 72 contained 0 incompatible loci -- more than two alleles present.
 25 Population 73 contained 0 incompatible loci -- more than two alleles present.
 26 Population 74 contained 0 incompatible loci -- more than two alleles present.
 27 Population 75 contained 0 incompatible loci -- more than two alleles present.
 28 Population 76 contained 0 incompatible loci -- more than two alleles present.
 29 Population 77 contained 0 incompatible loci -- more than two alleles present.
 30 Population 86 contained 0 incompatible loci -- more than two alleles present.
 31 Population 87 contained 0 incompatible loci -- more than two alleles present.
 32 Population 88 contained 0 incompatible loci -- more than two alleles present.
 33 Population 89 contained 0 incompatible loci -- more than two alleles present.
 34 Population 90 contained 0 incompatible loci -- more than two alleles present.
 35 Population 91 contained 0 incompatible loci -- more than two alleles present.
 36 Population 92 contained 0 incompatible loci -- more than two alleles present.
 37 Population 93 contained 0 incompatible loci -- more than two alleles present.
 38
 39 # Distribution of the number of SNPs per locus.
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60

#n_snps	n_loci
0	1774
1	3836
2	2916
3	2067
4	1465
5	1004
6	696
7	424
8	323

1		
2		
3	9	231
4	10	106
5	11	87
6	12	61
7	13	37
8	14	28
9	15	26
10	16	18
11	17	5
12	18	6
13	19	4
14	20	5
15	21	2
16	22	1
17	23	1
18	24	1
19	25	0
20	26	0
21	27	1
22	28	0
23	29	1
24		
25		
26		
27		
28		
29		
30		

Populations log for dataset 2

Distribution of population loci.

Distribution of valid loci matched to catalog locus.

Valid samples at locus Count

36	1	313602
37	2	65350
38	3	18307
39	4	8690
40	5	5498
41	6	4104
42	7	3142
43	8	2794
44	9	2283
45	10	2055
46	11	1949
47	12	1761
48	13	1652
49	14	1547
50	15	1455
51	16	1439
52	17	1342
53		
54		
55		
56		
57		
58		
59		
60		

1		
2		
3	18	1296
4	19	1283
5	20	1190
6	21	1229
7	22	1175
8	23	1156
9	24	1093
10	25	1088
11	26	1094
12	27	1030
13	28	1003
14	29	1016
15	30	960
16	31	967
17	32	918
18	33	950
19	34	990
20	35	977
21	36	933
22	37	830
23	38	859
24	39	910
25	40	898
26	41	838
27	42	848
28	43	830
29	44	852
30	45	869
31	46	868
32	47	804
33	48	817
34	49	888
35	50	827
36	51	849
37	52	880
38	53	778
39	54	812
40	55	891
41	56	838
42	57	873
43	58	840
44	59	831
45	60	906
46	61	947
47	62	958
48	63	952
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		
60		

For Review Only

1
2
3 64 913
4 65 1069
5 66 1128
6 67 1178
7 68 1243
8 69 1275
9 70 1335
10 71 1237
11 72 1163
12 73 998
13 74 626
14 75 400
15 76 326

18 # Distribution of confounded loci at catalog locus.

19 # Confounded samples at locus Count

20 0 493502

21 # Distribution of missing loci at catalog loci.

22 # Absent samples at locus Count

23 0 493502

24 # Distribution of population loci after applying locus constraints.

25 # Distribution of valid loci matched to catalog locus.

26 # Valid samples at locus Count

27 45 671

28 46 710

29 47 735

30 48 703

31 49 681

32 50 721

33 51 713

34 52 640

35 53 677

36 54 681

37 55 737

38 56 711

39 57 722

40 58 721

41 59 768

42 60 858

43 61 750

44 62 765

45 63 691

46 64 699

47 65 667

48 66 582

49 67 516

50 68 433

51
52
53
54
55
56
57
58
59
60

1		
2		
3	69	305
4	70	241
5	71	171
6	72	140
7	73	120
8	74	115
9	75	151
10	76	145

Distribution of confounded loci at catalog locus.

# Confounded samples at locus	Count
-------------------------------	-------

0	17940
---	-------

Distribution of missing loci at catalog loci.

# Absent samples at locus	Count
---------------------------	-------

0	309
1	282
2	352
3	547
4	847
5	1123
6	1390
7	1588
8	1603
9	1763
10	1710
11	1501
12	1278
13	1122
14	858
15	598
16	425
17	289
18	173
19	96
20	51
21	21
22	9
23	4
25	1

Population 1 contained 0 incompatible loci -- more than two alleles present.

Population 2 contained 0 incompatible loci -- more than two alleles present.

Population 3 contained 0 incompatible loci -- more than two alleles present.

Population 4 contained 0 incompatible loci -- more than two alleles present.

Population 5 contained 0 incompatible loci -- more than two alleles present.

Population 6 contained 0 incompatible loci -- more than two alleles present.

Population 7 contained 0 incompatible loci -- more than two alleles present.

Population 8 contained 0 incompatible loci -- more than two alleles present.

56
57
58
59
60

1
2
3 Population 9 contained 0 incompatible loci -- more than two alleles present.
4 Population 10 contained 0 incompatible loci -- more than two alleles present.
5 Population 11 contained 0 incompatible loci -- more than two alleles present.
6 Population 12 contained 0 incompatible loci -- more than two alleles present.
7 Population 13 contained 0 incompatible loci -- more than two alleles present.
8 Population 14 contained 0 incompatible loci -- more than two alleles present.
9 Population 15 contained 0 incompatible loci -- more than two alleles present.
10 Population 16 contained 0 incompatible loci -- more than two alleles present.
11 Population 17 contained 0 incompatible loci -- more than two alleles present.
12 Population 18 contained 0 incompatible loci -- more than two alleles present.
13 Population 19 contained 0 incompatible loci -- more than two alleles present.
14 Population 20 contained 0 incompatible loci -- more than two alleles present.
15 Population 21 contained 0 incompatible loci -- more than two alleles present.
16 Population 22 contained 0 incompatible loci -- more than two alleles present.
17 Population 25 contained 0 incompatible loci -- more than two alleles present.
18 Population 26 contained 0 incompatible loci -- more than two alleles present.
19 Population 30 contained 0 incompatible loci -- more than two alleles present.
20 Population 31 contained 0 incompatible loci -- more than two alleles present.
21 Population 32 contained 0 incompatible loci -- more than two alleles present.
22 Population 33 contained 0 incompatible loci -- more than two alleles present.
23 Population 34 contained 0 incompatible loci -- more than two alleles present.
24 Population 35 contained 0 incompatible loci -- more than two alleles present.
25 Population 36 contained 0 incompatible loci -- more than two alleles present.
26 Population 37 contained 0 incompatible loci -- more than two alleles present.
27 Population 38 contained 0 incompatible loci -- more than two alleles present.
28 Population 39 contained 0 incompatible loci -- more than two alleles present.
29 Population 40 contained 0 incompatible loci -- more than two alleles present.
30 Population 41 contained 0 incompatible loci -- more than two alleles present.
31 Population 45 contained 0 incompatible loci -- more than two alleles present.
32 Population 46 contained 0 incompatible loci -- more than two alleles present.
33 Population 47 contained 0 incompatible loci -- more than two alleles present.
34 Population 48 contained 0 incompatible loci -- more than two alleles present.
35 Population 49 contained 0 incompatible loci -- more than two alleles present.
36 Population 50 contained 0 incompatible loci -- more than two alleles present.
37 Population 51 contained 0 incompatible loci -- more than two alleles present.
38 Population 52 contained 0 incompatible loci -- more than two alleles present.
39 Population 53 contained 0 incompatible loci -- more than two alleles present.
40 Population 54 contained 0 incompatible loci -- more than two alleles present.
41 Population 55 contained 0 incompatible loci -- more than two alleles present.
42 Population 56 contained 0 incompatible loci -- more than two alleles present.
43 Population 57 contained 0 incompatible loci -- more than two alleles present.
44 Population 58 contained 0 incompatible loci -- more than two alleles present.
45 Population 59 contained 0 incompatible loci -- more than two alleles present.
46 Population 60 contained 0 incompatible loci -- more than two alleles present.
47 Population 61 contained 0 incompatible loci -- more than two alleles present.
48 Population 62 contained 0 incompatible loci -- more than two alleles present.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Population 63 contained 0 incompatible loci -- more than two alleles present.
4 Population 64 contained 0 incompatible loci -- more than two alleles present.
5 Population 65 contained 0 incompatible loci -- more than two alleles present.
6 Population 66 contained 0 incompatible loci -- more than two alleles present.
7 Population 68 contained 0 incompatible loci -- more than two alleles present.
8 Population 69 contained 0 incompatible loci -- more than two alleles present.
9 Population 70 contained 0 incompatible loci -- more than two alleles present.
10 Population 71 contained 0 incompatible loci -- more than two alleles present.
11 Population 72 contained 0 incompatible loci -- more than two alleles present.
12 Population 73 contained 0 incompatible loci -- more than two alleles present.
13 Population 74 contained 0 incompatible loci -- more than two alleles present.
14 Population 75 contained 0 incompatible loci -- more than two alleles present.
15 Population 76 contained 0 incompatible loci -- more than two alleles present.
16 Population 77 contained 0 incompatible loci -- more than two alleles present.
17 Population 86 contained 0 incompatible loci -- more than two alleles present.
18 Population 87 contained 0 incompatible loci -- more than two alleles present.
19 Population 88 contained 0 incompatible loci -- more than two alleles present.
20 Population 89 contained 0 incompatible loci -- more than two alleles present.
21 Population 90 contained 0 incompatible loci -- more than two alleles present.
22 Population 91 contained 0 incompatible loci -- more than two alleles present.
23 Population 92 contained 0 incompatible loci -- more than two alleles present.
24 Population 93 contained 0 incompatible loci -- more than two alleles present.
25 Population 1 contained 0 incompatible loci -- more than two alleles present.
26 Population 2 contained 0 incompatible loci -- more than two alleles present.
27 Population 3 contained 0 incompatible loci -- more than two alleles present.
28 Population 4 contained 0 incompatible loci -- more than two alleles present.
29 Population 5 contained 0 incompatible loci -- more than two alleles present.
30 Population 6 contained 0 incompatible loci -- more than two alleles present.
31 Population 7 contained 0 incompatible loci -- more than two alleles present.
32 Population 8 contained 0 incompatible loci -- more than two alleles present.
33 Population 9 contained 0 incompatible loci -- more than two alleles present.
34 Population 10 contained 0 incompatible loci -- more than two alleles present.
35 Population 11 contained 0 incompatible loci -- more than two alleles present.
36 Population 12 contained 0 incompatible loci -- more than two alleles present.
37 Population 13 contained 0 incompatible loci -- more than two alleles present.
38 Population 14 contained 0 incompatible loci -- more than two alleles present.
39 Population 15 contained 0 incompatible loci -- more than two alleles present.
40 Population 16 contained 0 incompatible loci -- more than two alleles present.
41 Population 17 contained 0 incompatible loci -- more than two alleles present.
42 Population 18 contained 0 incompatible loci -- more than two alleles present.
43 Population 19 contained 0 incompatible loci -- more than two alleles present.
44 Population 20 contained 0 incompatible loci -- more than two alleles present.
45 Population 21 contained 0 incompatible loci -- more than two alleles present.
46 Population 22 contained 0 incompatible loci -- more than two alleles present.
47 Population 25 contained 0 incompatible loci -- more than two alleles present.
48 Population 26 contained 0 incompatible loci -- more than two alleles present.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Population 88 contained 0 incompatible loci -- more than two alleles present.
4 Population 89 contained 0 incompatible loci -- more than two alleles present.
5 Population 90 contained 0 incompatible loci -- more than two alleles present.
6 Population 91 contained 0 incompatible loci -- more than two alleles present.
7 Population 92 contained 0 incompatible loci -- more than two alleles present.
8 Population 93 contained 0 incompatible loci -- more than two alleles present.

9 # Distribution of the number of SNPs per locus.

10 #n_snps n_loci
11
12 0 1774
13 1 13352
14

15
16 ##### Populations log for dataset 3 #####
17

18 # Note: Individual distributions can be extracted using the `stacks-dist-extract` utility.

19 # e.g. `stacks-dist-extract populations.log.distrib dist_name`
20

21 BEGIN batch_progress

22 C01: analyzed 19835 loci; filtered 230810 loci; 250645 loci seen.

23 1931548 genomic sites, of which 56414 were covered by multiple loci (2.9%).

24 C02: analyzed 22724 loci; filtered 259151 loci; 281875 loci seen.

25 2206248 genomic sites, of which 70104 were covered by multiple loci (3.2%).

26 C03: analyzed 30062 loci; filtered 351403 loci; 381465 loci seen.

27 2926008 genomic sites, of which 85985 were covered by multiple loci (2.9%).

28 C04: analyzed 21664 loci; filtered 249411 loci; 271075 loci seen.

29 2092935 genomic sites, of which 74845 were covered by multiple loci (3.6%).

30 C05: analyzed 17250 loci; filtered 191665 loci; 208915 loci seen.

31 1672712 genomic sites, of which 53943 were covered by multiple loci (3.2%).

32 C06: analyzed 21088 loci; filtered 240522 loci; 261610 loci seen.

33 2046130 genomic sites, of which 65484 were covered by multiple loci (3.2%).

34 C07: analyzed 25390 loci; filtered 288384 loci; 313774 loci seen.

35 2470551 genomic sites, of which 71939 were covered by multiple loci (2.9%).

36 C08: analyzed 21791 loci; filtered 247627 loci; 269418 loci seen.

37 2124977 genomic sites, of which 57141 were covered by multiple loci (2.7%).

38 C09: analyzed 20756 loci; filtered 240403 loci; 261159 loci seen.

39 2015736 genomic sites, of which 63168 were covered by multiple loci (3.1%).

40 END batch_progress
41
42
43
44

45 BEGIN samples_per_loc_prefilters

46 # Distribution of valid samples matched to a catalog locus prior to filtering.

47 n_samples n_loci
48
49 1 419184
50 2 231523
51 3 196563
52 4 177894
53 5 161952
54 6 145719
55
56
57
58
59
60

1		
2		
3	7	126487
4	8	111201
5	9	98979
6	10	87384
7	11	78741
8	12	72739
9	13	69163
10	14	62672
11	15	56675
12	16	50546
13	17	44467
14	18	38520
15	19	32877
16	20	27703
17	21	23793
18	22	20021
19	23	16651
20	24	14078
21	25	11455
22	26	9373
23	27	7816
24	28	6471
25	29	5255
26	30	4498
27	31	4002
28	32	3471
29	33	3163
30	34	2762
31	35	2479
32	36	2346
33	37	2106
34	38	2052
35	39	1798
36	40	1692
37	41	1650
38	42	1584
39	43	1532
40	44	1487
41	45	1428
42	46	1462
43	47	1315
44	48	1316
45	49	1268
46	50	1234
47	51	1253
48	52	1210
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		
60		

For Review Only

```
1
2
3      53      1157
4      54      1271
5      55      1138
6      56      1250
7      57      1262
8      58      1305
9      59      1267
10     60      1212
11     61      1265
12     62      1317
13     63      1342
14     64      1374
15     65      1333
16     66      1440
17     67      1484
18     68      1489
19     69      1623
20     70      1801
21     71      2015
22     72      2202
23     73      2785
24     74      3398
25     75      5041
26     76      7155
27
28
29
30
31 END samples_per_loc_prefilters
32
33 BEGIN missing_samples_per_loc_prefilters
34 # Distribution of missing samples for each catalog locus prior to filtering.
35 # Absent samples at locus      Count
36 0      7155
37 1      5041
38 2      3398
39 3      2785
40 4      2202
41 5      2015
42 6      1801
43 7      1623
44 8      1489
45 9      1484
46 10     1440
47 11     1333
48 12     1374
49 13     1342
50 14     1317
51 15     1265
52 16     1212
53
54
55
56
57
58
59
60
```

For Review Only

1		
2		
3	17	1267
4	18	1305
5	19	1262
6	20	1250
7	21	1138
8	22	1271
9	23	1157
10	24	1210
11	25	1253
12	26	1234
13	27	1268
14	28	1316
15	29	1315
16	30	1462
17	31	1428
18	32	1487
19	33	1532
20	34	1584
21	35	1650
22	36	1692
23	37	1798
24	38	2052
25	39	2106
26	40	2346
27	41	2479
28	42	2762
29	43	3163
30	44	3471
31	45	4002
32	46	4498
33	47	5255
34	48	6471
35	49	7816
36	50	9373
37	51	11455
38	52	14078
39	53	16651
40	54	20021
41	55	23793
42	56	27703
43	57	32877
44	58	38520
45	59	44467
46	60	50546
47	61	56675
48	62	62672
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		
60		

For Review Only

```
1
2
3      63      69163
4      64      72739
5      65      78741
6      66      87384
7      67      98979
8      68     111201
9      69     126487
10     70     145719
11     71     161952
12     72     177894
13     73     196563
14     74     231523
15     75     419184
```

```
18 END missing_samples_per_loc_prefilters
```

```
20 BEGIN snps_per_loc_prefilters
```

```
21 # Distribution of the number of SNPs per catalog locus prior to filtering.
```

```
22 n_snps n_loci
23 0      2036813
24 1      220048
25 2      102200
26 3      54511
27 4      31646
28 5      19251
29 6      12286
30 7      8067
31 8      5338
32 9      3388
33 10     2338
34 11     1556
35 12     998
36 13     565
37 14     358
38 15     239
39 16     120
40 17     78
41 18     49
42 19     30
43 20     21
44 21     12
45 22     8
46 23     5
47 24     3
48 25     4
49 26     2
50 28     1
```

```
56
57
58
59
60
```

```
1
2
3      32      1
4      END snps_per_loc_prefilters
5
6      BEGIN samples_per_loc_postfilters
7      # Distribution of valid samples matched to a catalog locus after filtering.
8      n_samples      n_loci
9
10     2      13122
11     3      12184
12     4      78240
13     5      308
14     6      4485
15     7      3006
16     8      16051
17     9      325
18    10      2454
19    11      1829
20    12      5359
21    13      407
22    14      1909
23    15      1474
24    16      2433
25    17      558
26    18      1533
27    19      1231
28    20      1373
29    21      600
30    22      1329
31    23      1092
32    24      948
33    25      703
34    26      1197
35    27      927
36    28      675
37    29      810
38    30      1036
39    31      745
40    32      516
41    33      888
42    34      932
43    35      647
44    36      434
45    37      967
46    38      921
47    39      520
48    40      447
49    41      1051
50
51
52
53
54
55
56
57
58
59
60
```

1
2
3 42 924
4 43 359
5 44 524
6 45 1159
7 46 720
8 47 244
9 48 661
10 49 1161
11 50 589
12 51 173
13 52 828
14 53 1174
15 54 534
16 55 114
17 56 1188
18 57 1123
19 58 375
20 59 65
21 60 1580
22 61 1060
23 62 249
24 63 33
25 64 2218
26 65 957
27 66 136
28 67 17
29 68 3278
30 69 719
31 70 72
32 71 3
33 72 5036
34 73 457
35 74 9
36 76 7155

42 END samples_per_loc_postfilters

44 BEGIN missing_samples_per_loc_postfilters

45 # Distribution of missing samples for each catalog locus after filtering.

46 # Absent samples at locus Count

47 0 7155
48 2 9
49 3 457
50 4 5036
51 5 3
52 6 72
53 7 719

54
55
56
57
58
59
60

For Review Only

1		
2		
3	8	3278
4	9	17
5	10	136
6	11	957
7	12	2218
8	13	33
9	14	249
10	15	1060
11	16	1580
12	17	65
13	18	375
14	19	1123
15	20	1188
16	21	114
17	22	534
18	23	1174
19	24	828
20	25	173
21	26	589
22	27	1161
23	28	661
24	29	244
25	30	720
26	31	1159
27	32	524
28	33	359
29	34	924
30	35	1051
31	36	447
32	37	520
33	38	921
34	39	967
35	40	434
36	41	647
37	42	932
38	43	888
39	44	516
40	45	745
41	46	1036
42	47	810
43	48	675
44	49	927
45	50	1197
46	51	703
47	52	948
48	53	1092
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		
60		

For Review Only

```
1
2
3      54      1329
4      55      600
5      56      1373
6      57      1231
7      58      1533
8      59      558
9      60      2433
10     61      1474
11     62      1909
12     63      407
13     64      5359
14     65      1829
15     66      2454
16     67      325
17     68      16051
18     69      3006
19     70      4485
20     71      308
21     72      78240
22     73      12184
23     74      13122
```

```
24
25
26 END missing_samples_per_loc_postfilters
```

```
27
28
29 BEGIN snps_per_loc_postfilters
```

```
30 # Distribution of the number of SNPs per catalog locus (after filtering).
```

```
31 n_snps n_loci
32
33 0      147021
34 1      21457
35 2      11928
36 3      7243
37 4      4795
38 5      3005
39 6      1886
40 7      1236
41 8      765
42 9      498
43 10     278
44 11     179
45 12     111
46 13     73
47 14     41
48 15     21
49 16     11
50 17     3
51 18     3
52 19     6
```

```
53
54
55
56
57
58
59
60
```

1
2
3 END snps_per_loc_postfilters
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only

1
2
3
4 **Title: Supplementary Information for: Feral populations of *Brassica oleracea* along Atlantic**
5
6 **coasts in western Europe**
7
8
9

10
11 **Table 1:** Information about *Brassica oleracea* plants from the Atlantic coasts in western Europe from
12 double-digest restriction-site associated DNA sequencing, including: the number of reads obtained from
13 sequencing, the number of reads that mapped to the *B. oleracea* genome and the percentage of reads
14 that mapped to the *B. oleracea* genome.

Region	Population	Number of reads	Number of reads mapped	Percentage of reads mapped
East Scotland	Auchmithie	1387753	1207726	87.03
East Scotland	Auchmithie	1567016	1357070	86.60
East Scotland	Auchmithie	1464345	1268561	86.63
East Scotland	Auchmithie	1621739	1403811	86.56
East Scotland	Crail	1595187	1378333	86.41
East Scotland	Crail	1468764	1288648	87.74
East Scotland	Crail	1592711	1400314	87.92
East Scotland	Crail	1427818	1261920	88.38
North-east Scotland	Fortrose	1793049	1492860	83.26
North-east Scotland	Fortrose	745	566	75.97
North-east Scotland	Fortrose	1418162	1169961	82.50
North-east Scotland	Fortrose	1015983	871504	85.78
West Scotland	Kildonan Castle	1920023	1675933	87.29
West Scotland	Kildonan Castle	827008	725225	87.69
West Scotland	Kildonan Castle	863404	753105	87.23
West Scotland	Kildonan Castle	384814	334230	86.85
North-east England	Tynemouth	5361799	4611195	86.00
North-east England	Tynemouth	730526	627536	85.90
North-east England	Tynemouth	1345547	1151687	85.59
North-east England	Tynemouth	2928179	2492162	85.11
North-east England	Staithe	1147	953	83.09
North-east England	Staithe	873	718	82.25
North-east England	Staithe	1443	1247	86.42
North-east England	Staithe	910	752	82.64
North-east England	Whitby	1126643	972813	86.35
North-east England	Whitby	2268574	1943089	85.65
North-east England	Whitby	1119790	977734	87.31

Table 1: Information about sequencing of *Brassica oleracea* plants cont.

Region	Population	Number of reads	Number of reads mapped	Percentage of reads mapped
North-east England	Whitby	1063470	905937	85.19
North Wales	Little Orme	410	332	80.98
North Wales	Little Orme	424	356	83.96
North Wales	Little Orme	470	389	82.77
North Wales	Little Orme	601	479	79.70
South Wales	Tenby	1462901	1245452	85.14
South Wales	Tenby	2222716	1882245	84.68
South Wales	Tenby	1593699	1340950	84.14
South Wales	Tenby	1282352	1083760	84.51
South Wales	Llantwit Major	1329946	1146093	86.18
South Wales	Llantwit Major	1396788	1210897	86.69
South Wales	Llantwit Major	1919399	1650157	85.97
South Wales	Llantwit Major	2591684	2222948	85.77
South-west England	Prussia cove	1619523	1387691	85.69
South-west England	Prussia cove	1772075	1513752	85.42
South-west England	Prussia cove	1306856	1127550	86.28
South-west England	Prussia cove	930252	801321	86.14
South-west England	Fowey	1202859	1043323	86.74
South-west England	Fowey	2800428	2420953	86.45
South-west England	Fowey	1096937	951131	86.71
South-west England	Fowey	655043	568748	86.83
South-west England	West Looe	526143	456982	86.86
South-west England	West Looe	1275394	1099013	86.17
South-west England	St Aldeham's Head	1996066	1733925	86.87
South-west England	St Aldeham's Head	1826145	1597182	87.46
South-west England	St Aldeham's Head	1373774	1199240	87.30
South-west England	St Aldeham's Head	1735302	1506195	86.80
Asturias	Cudillero	298891	263713	88.23
Asturias	Cudillero	1296387	1139395	87.89
Asturias	Cudillero	947934	835735	88.16
Asturias	Cudillero	3575040	3158715	88.35
Asturias	Playa de Xágo	1812681	1573772	86.82
Asturias	Playa de Xágo	1162553	1003317	86.30

Table 1: Information about sequencing of *Brassica oleracea* plants cont.

Region	Population	Number of reads	Number of reads mapped	Percentage of reads mapped
Asturias	Playa de Xágo	578767	495639	85.64
Asturias	Playa de Xágo	1278240	1082892	84.72
Asturias	Cabo de Peñas	1028128	893090	86.87
Asturias	Cabo de Peñas	1268678	1082179	85.30
Asturias	Cabo de Peñas	2897986	2476276	85.45
Asturias	Cabo de Peñas	413511	359603	86.96
Asturias	Playas de Viodo	482032	422063	87.56
Asturias	Playas de Viodo	737102	652759	88.56
Asturias	Playas de Viodo	220363	192283	87.26
Asturias	Playas de Viodo	721625	624490	86.54
Asturias	Tazonos	710	601	84.65
Asturias	Tazonos	1922	1587	82.57
Asturias	Tazonos	1070	918	85.79
Cantabrica	Playa Pedrero	1009564	867686	85.95
Cantabrica	Playa Pedrero	1298664	1129231	86.95
Cantabrica	Playa Pedrero	3058101	2633809	86.13
Cantabrica	Playa Pedrero	1414237	1221669	86.38
Cantabrica	La Franca	604	509	84.27
Cantabrica	La Franca	466	379	81.33
Cantabrica	La Franca	520	433	83.27
Basque	San Juan de Gaxtelugatxe	1444693	1253303	86.75
Basque	San Juan de Gaxtelugatxe	3135997	2723927	86.86
Basque	San Juan de Gaxtelugatxe	1527228	1315084	86.11
Basque	San Juan de Gaxtelugatxe	2534715	2178782	85.96
Basque	Getaría	1419546	1232331	86.81
Basque	Getaría	1146516	994255	86.72
Basque	Getaría	1325400	1133457	85.52
Basque	Getaría	1304818	1117297	85.63
Basque	San Sebastian	1203388	1051422	87.37
Basque	San Sebastian	2227159	1921417	86.27
Basque	San Sebastian	898	612	68.15
Basque	San Sebastian	964	798	82.78

Table 1: Information about sequencing of *Brassica oleracea* plants cont.

Region	Population	Number of reads	Number of reads mapped	Percentage of reads mapped
Basque	San Sebastian	2768379	2426913	87.67

Table 2: Number and percentage of missing SNPs for each population in dataset 1, which contains 42,517 SNPs.

Region	Population	Number of SNPs	Number of missing SNPs	Percentage of missing SNPs
East Scotland	Auchmithie	35284	7233	17.0
East Scotland	Auchmithie	34909	7608	17.9
East Scotland	Auchmithie	35374	7143	16.8
East Scotland	Auchmithie	35936	6581	15.5
East Scotland	Crail	37553	4964	11.7
East Scotland	Crail	37964	4553	10.7
East Scotland	Crail	36184	6333	14.9
East Scotland	Crail	35580	6937	16.3
North-east Scotland	Fortrose	31267	11250	26.5
North-east Scotland	Fortrose	24585	17932	42.2
North-east Scotland	Fortrose	24019	18498	43.5
West Scotland	Kildonan Castle	37569	4948	11.6
West Scotland	Kildonan Castle	25978	16539	38.9
West Scotland	Kildonan Castle	24633	17884	42.1
West Scotland	Kildonan Castle	8667	33850	79.6
North-east England	Tynemouth	39926	2591	6.1
North-east England	Tynemouth	11336	31181	73.3
North-east England	Tynemouth	23338	19179	45.1
North-east England	Tynemouth	35321	7196	16.9
North-east England	Whitby	32574	9943	23.4
North-east England	Whitby	38987	3530	8.3
North-east England	Whitby	33152	9365	22.0
North-east England	Whitby	30207	12310	29.0
South Wales	Tenby	27058	15459	36.4
South Wales	Tenby	30274	12243	28.8
South Wales	Tenby	22793	19724	46.4
South Wales	Tenby	25416	17101	40.2

Table 2: Information about dataset 1 SNPs cont.

Region	Population	Number of SNPs	Number of missing SNPs	Percentage of missing SNPs
South Wales	Llantwit Major	34529	7988	18.8
South Wales	Llantwit Major	35783	6734	15.8
South Wales	Llantwit Major	37567	4950	11.6
South Wales	Llantwit Major	38861	3656	8.6
South-west England	Prussia cove	35111	7406	17.4
South-west England	Prussia cove	33978	8539	20.1
South-west England	Prussia cove	30779	11738	27.6
South-west England	Prussia cove	17982	24535	57.7
South-west England	Fowey	28968	13549	31.9
South-west England	Fowey	39332	3185	7.5
South-west England	Fowey	28774	13743	32.3
South-west England	Fowey	16316	26201	61.6
South-west England	West Looe	25330	17187	40.4
South-west England	West Looe	33283	9234	21.7
South-west England	St Aldeham's Head	38935	3582	8.4
South-west England	St Aldeham's Head	38055	4462	10.5
South-west England	St Aldeham's Head	36817	5700	13.4
South-west England	St Aldeham's Head	36176	6341	14.9
Asturias	Cudillero	11931	30586	71.9
Asturias	Cudillero	35501	7016	16.5
Asturias	Cudillero	3294	9593	22.6
Asturias	Cudillero	39800	2717	6.4
Asturias	Playa de Xágo	37756	4761	11.2
Asturias	Playa de Xágo	33977	8540	20.1
Asturias	Playa de Xágo	19052	23465	55.2
Asturias	Playa de Xágo	32265	10252	24.1
Asturias	Cabo de Peñas	33211	9306	21.9
Asturias	Cabo de Peñas	29122	13395	31.5
Asturias	Cabo de Peñas	38730	3787	8.9
Asturias	Cabo de Peñas	8223	34294	80.7
Asturias	Playas de Viodo	11282	31235	73.5
Asturias	Playas de Viodo	26662	15855	37.3
Asturias	Playas de Viodo	1731	40786	95.9

Table 2: Information about dataset 1 SNPs cont.

Region	Population	Number of SNPs	Number of missing SNPs	Percentage of missing SNPs
Asturias	Playas de Viedo	19420	23097	54.3
Cantabrica	Playa Pedrero	31359	11158	26.2
Cantabrica	Playa Pedrero	35116	7401	17.4
Cantabrica	Playa Pedrero	39848	2669	6.3
Cantabrica	Playa Pedrero	36310	6207	14.6
Basque	San Juan de Gaxtelugatxe	34359	8158	19.2
Basque	San Juan de Gaxtelugatxe	39746	2771	6.5
Basque	San Juan de Gaxtelugatxe	33073	9444	22.2
Basque	San Juan de Gaxtelugatxe	37287	5230	12.3
Basque	Getaría	34975	7542	17.7
Basque	Getaría	28516	14001	32.9
Basque	Getaría	31536	10981	25.8
Basque	Getaría	33376	9141	21.5
Basque	San Sebastian	34737	7780	18.3
Basque	San Sebastian	39028	3489	8.2
Basque	San Sebastian	39569	2948	6.9

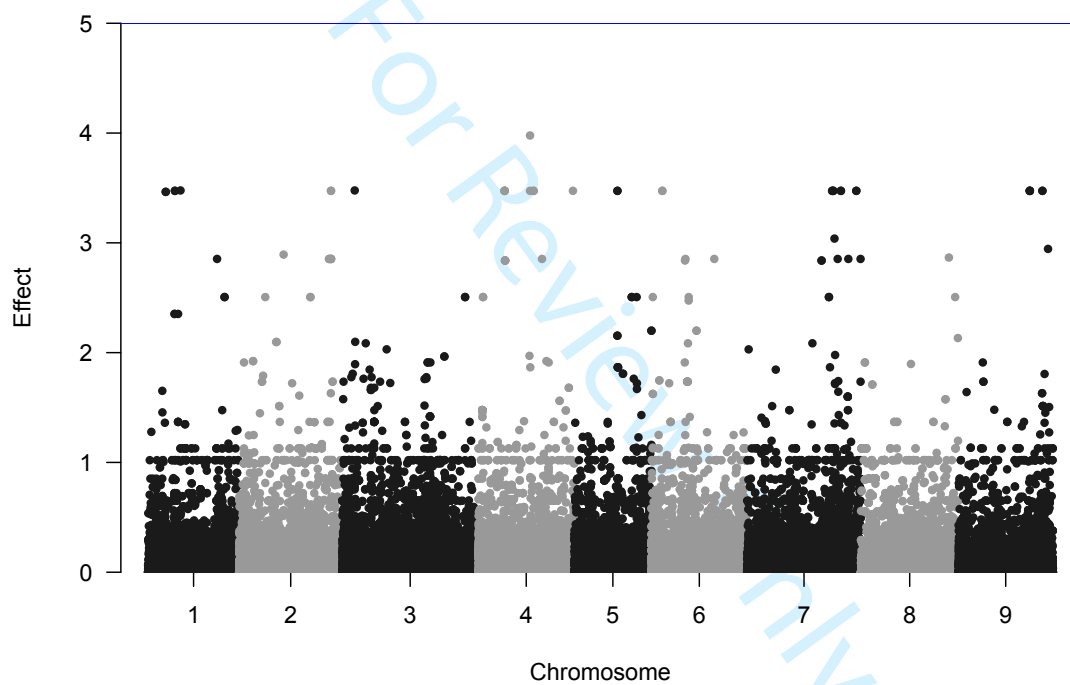


Figure 1: Manhattan plot showing the distribution of SNPs from ddRAD sequencing of *B. oleracea* plants from natural populations in western Europe mapped to the *B. oleracea* genome from Liu *et al.* 2014.

Mittell et al. Response to reviewers July 2020

Reviewer Comments to Author:

Reviewer: 1

Comments to the Author

In this ms Mittell et al examine the population structure of *Brassica oleracea* putatively wild populations from the Atlantic range of the species, with the aim of discerning their wild or feral origin, as well as to identify signatures of local adaptation. For this, they obtained ca. 13,000 ddRAD loci on 76 individuals of 20 populations of UK and Spain. The ms is nicely written and tells a congruent story. The molecular and bioinformatic methods and data analyses are also well described and employed. My only major concern is that the sampling design is not adequate to test the main driving question of the ms, as I explain further below.

MAJOR

The main question of the ms is what is the mostly likely origin (wild or feral) of wild *B. oleracea* in the Atlantic region. To examine this, the authors perform phylogenetic and population structure analyses with a comprehensive sampling of these potentially wild populations, even incorporating data on the first local records of the plants. This is all very good and informative, and I agree with the way the data was analysed and discussed under the umbrella of the driving question. However, in my opinion results like the weird no-isolation-by-distance could also be explained by an intricate phylogeographic history, and not necessarily by human movement of populations. Therefore, I think it would be much more robust to include cultivar samples (specially kale as mentioned in l. 362) and an outgroup in the analyses. Cultivated samples would be expected to be a subset of the wild variation, and feral populations (despite how "wild" they look phenotypically) would be expected to be a subset of the cultivated variation, as has been shown in other wild-domesticated species complexes.

I know it may be unrealistic for the work group to sample and sequence again, so I dare not ask for that. But given that this is a well studied crop, there is already available data on cultivated forms (e.g. <https://www.nature.com/articles/s41438-018-0040-3>) that with some bioinformatics wrangling could be made comparable with the present dataset, at least to obtain enough loci to perform a phylogenetic analyses.

If, for whatever reason, incorporating published data from domesticated forms is not possible, I think the authors should at least acknowledge this caveat and re-focus the paper to make the wild-feral point less central. I'm sure this is possible, given that the dataset is very good and relevant to CWR research.

Reply: Pages 6 and 15. Thank you for your comments here. We have downplayed focus on the question of origins by: (1) being more cautious about what we can conclude in the discussion; and (2) altering our objectives to focus more explicitly on investigating genetic variation and population structure within the two focal regions sampled. We acknowledge the lack of cultivars in the discussion (lines 384--387) by saying explicitly that it would be interesting to add to our work by identifying, sequencing and comparing results with local cultivars. There were some cultivars included in the original sequencing. However, when some of the sequencing failed (the reason for not all the natural localities being included), the cultivar samples were unfortunately amongst these. Since our main focus was on comparing the wild populations in the UK and Spain, we had maximised the number of populations and samples that could be included from those. In the future we agree that it would be useful to include many more cultivars to cover this aspect comprehensively and avoid issues with loss of samples during sequencing.

1
2
3
4 Carrying out another study with the additional analyses suggested would be a good place to start for
5 future work. Using published data is an intriguing idea but combining genotype by sequencing data
6 generated using different restriction enzymes would not be trivial. Instead, we have refocused the
7 aims of the paper to emphasise adding knowledge about diversity in the wild.
8
9

10 DATA ACCESSIBILITY

11
12 Please notice that raw data should be deposited at SRA (not Dryad), and that your Dryad repository
13 should contain Stacks genotypes output (along with any other relevant to downstream analyses) and
14 scripts (bash scripts for processing data and R scripts for the mentioned analyses).

15 Reply: Page 22. Thank you, we have deposited the raw sequencing data and associated meta data on
16 the European Nucleotide Archive under the study accession number: PRJEB38464
17 ([\url{http://www.ebi.ac.uk/ena/data/view/PRJEB38464}](http://www.ebi.ac.uk/ena/data/view/PRJEB38464)).
18
19

20 MINOR

21
22 * The introduction, specially the first paragraph, is a bit too focused on the genetic consequences of
23 domestication due to artificial selection, and it misses to explicitly incorporate the effect of genetic
24 drift.

25 Reply: Page 2. Increased genetic drift is included as an example of what can happen due to artificial
26 selection and domestication bottlenecks in lines 46--48.
27
28

29 * Last paragraph of the introduction, main question addressed: I suggest changing "(2) how much
30 population structure exists between geographically isolated populations" for "(2) what is the
31 population structure among the Atlantic range of the species and how much differentiation exists
32 between isolated populations?"

33 Reply: Page 6. We have changed our objectives to "The following questions were addressed: (1)
34 how much genetic variation exists among wild populations of *B. oleracea* in the UK and Spain; (2)
35 how are populations structured in the Atlantic region and how much differentiation exists between
36 isolated populations? and (3) are there signals of local adaptation to the environment?" Lines 143--
37 145
38
39

40 * It would be useful to incorporate few references to other wild-domesticated species complexes
41 where the origin of wild/feral populations was disentangled using genomic data.

42 Reply: Page 4. We have added the following into the introduction: "Compared to domestication,
43 feralization is under-investigated; however modern genomic data are allowing its occurrence to be
44 identified and consequences better understood (see examples in Henriksen 2018)." Line 84--86
45
46

47 * L 338-353. This paragraph feels a bit disconnected from your results, specially the last sentence.
48 Please discuss the microsatellite and nDNA previous results in light of your ddRAD data. Also notice
49 that you could extract cpDNA from your dataset.

50 Reply: Page 15. We are sorry that we hadn't made the connection to our results clear enough.
51 Particularly since the previous chloroplast data were based on microsatellites, it would be difficult to
52 extract from the ddRAD data, which filter out repeats because of the problems of interpreting when
53 considering only 100 bp loci. Therefore, we have reduced this paragraph and included more specific
54 comparison with our results. Lines 388--399.

55 "The two populations with the C:04 haplotype were in Tyne & Wear, in the northeast of England; in
56 the current study, this area is represented by the Tynemouth and Whitby populations. In line with
57 the rarity of the chloroplast haplotypes identified in this region in the previous study, these two
58 populations clustered most closely with populations not sampled by Allender et al. (2007);..."
59
60

1
2
3
4 Thanks for the opportunity of review your work. I enjoyed it, and hope these comments help to
5 improve it.
6

7
8 Reviewer: 2
9

10 Comments to the Author

11 Mittell et al. report analyses of feral populations of *Brassica oleracea* using sub-genomic sequencing
12 (RAD-seq) paired-end reads generated from individuals distributed along the coasts of UK and Spain.
13 They concluded that the feral populations recently expanded from domesticated varieties and there
14 might be a possible connection between genetic diversity and local soil pH.
15

16 The following comments should be addressed.

17
18 1. The authors presented a background about domestication and resulting bottleneck and reduced
19 genetic diversity; on the other hand, they also indicated that *B. oleracea* includes a great degree
20 morphological diversity. The study here focused on feral populations, thus it is not clear how the
21 results can be used to address questions about genetic diversity in domesticated crop varieties. It
22 seems that there is a disconnect between the early parts of the introduction and the actual study of
23 this paper.

24
25 **Reply: We had not intended to imply that our data could be used to understand genetic diversity in**
26 **domesticated crops. However, we do think that the patterns of genetic diversity and population**
27 **structuring are intriguing in relation to understanding how much variation exists in “wild”**
28 **populations. We have revised the abstract, introduction and discussion to make it clearer that our**
29 **focus was on the wild populations, rather than understanding more about domesticated crops.**
30

31
32 2. The author also discussed different hypotheses of the origins of the feral populations in the
33 Atlantic regions in north-western Europe, and that the most plausible hypothesis is that the feral
34 populations were derived from domestic crop varieties. Therefore, the conclusion is largely
35 confirmatory, lacking sufficient originality.

36
37 **Reply: Page 3. Although there is evidence from other sources, there has been a lack of genetic**
38 **evidence for either hypothesis. We hope by adding the following to the end of this paragraph on line**
39 **73 in the introduction, that this has been made clearer: “However, the genetic status of *B. oleracea***
40 **in the Atlantic region is still an open question (*B. oleracea* is classified as a native species in the UK**
41 **and an alien species in Spain; Euro+Med PlantBase 2020).”**
42

43
44 3. The authors stated that dataset 2 contain one SNP per each RAD locus, to obtain unlinked loci.
45 However, they also reported that the number of SNPs is 13158 in dataset 2. It seems very puzzling
46 how over 13000 loci could be unlinked for a species with a genetic map of about 900 cM. Did the
47 authors mean something else when they said “unlinked”?

48
49 **Reply: Page 7. Thank you for pointing out the potential confusion here. This has been changed to:**
50 **“SNPs linked within each RAD locus were avoided by only retaining one SNP at random per locus;”**
51 **line 191**
52

53
54 4. If one hypothesis to be tested is that feral populations were derived from domesticated varieties,
55 the study should include a number of cultivated varieties, especially those that have been grown in
56 UK and Spain during the history of cultivation in these regions.

57
58 **Reply: We agree that carrying out the same sequencing on a range of cultivated varieties from the**
59 **local regions would enhance our work in the future. As described in the response to reviewer 1,**
60 **there were some cultivars included in the original sequencing. However, some of the sequencing**
failed (the reason for not all the natural localities being included), and cultivar samples were
unfortunately amongst these. We still consider the implications of our results for testing this

1
2
3 hypothesis but downplay what can be concluded without sequencing of a wide range of cultivars.
4 We have modified the aims and restructured the discussion to focus more clearly on genetic
5 diversity and population structure within the “wild” populations.
6

7
8 5. The analysis here used the reference *B. oleracea* genome to identify genetic variations among
9 individuals of feral populations. Could there be sequences from these wild individuals that are too
10 divergent from the reference genome to be mapped? This idea is at least in principle supported by
11 the observation that 11.6-31.8% of reads were not mapped to the *B. oleracea*. Is it possible that
12 some of these reads are highly divergent from the reference genome sequence?

13 **Reply:** This is a good point and as more reference genomes become available then it would be worth
14 re-running analyses such as ours. However, since the main focus of our study was to use ddRAD to
15 identify shared SNPs that could be used for population genetics analyses, excluding highly divergent
16 regions might actually be preferable (e.g., to avoid regions under divergent selection).
17

18
19 6. If such sequences exist, could some of the divergent sequences reveal additional population
20 structural information regarding these individuals? For example, could a subset of these individuals
21 exhibit great similarities or differences than the analyses here have shown? If yes, such additional
22 information could alter the conclusion.

23 **Reply:** This would be interesting to test but ddRAD sequencing might not be the best approach to
24 use for this type of analysis, due to the difficulty of testing for selection with only 100 bp and
25 filtering to consider only one SNP per RAD locus. The population genetics analyses used assume
26 neutrality and so the conservative approach of excluding highly divergent loci might actually help to
27 reduce the risk of interpreting patterns based on divergent selection in different populations.
28

29
30 7. It is possible that de novo assembly might reveal additional sequence variation than those
31 revealed mapping to the reference genome would have missed.

32 An all-against-all comparison of de novo assembled sequences might reveal additional population
33 structure information.

34 **Reply:** We agree that de novo assembly might increase the number of loci included, which would be
35 important for genome-wide association studies. However, the 13,352 SNPs used here (dataset 2) are
36 suitable for assessing the population structure. Previous work has shown that there is an effect on
37 population structure of increasing the number of SNPs at low numbers (e.g., <100 SNPs in Morin *et al.*
38 2009 *Mol. Ecol. Resources*: 9(1) 66-73; and 50 -- 3500 SNPs in Nazareno *et al.* *Mol. Ecol.*
39 *Resources*: 17(6) 1136-1147), however, the number of SNPs used here is well above the minimum
40 number of SNPs required for the analyses carried out.
41

42
43 8. It would also be informative to know the map positions of the reads from the analyzed individuals,
44 in a supplemental figure or table. For example, are the SNPs clustered in a small number of
45 chromosomal regions, and distributed throughout the genome?

46 **Reply:** A supplementary figure has been added (and referred to in the text) to show the distribution
47 of the SNPs from our study mapped to the reference *B. oleracea* genome. In addition, the output
48 from stacks that contains the locations of the SNPs has been included in the supplementary material.
49

50
51 9. The RAD-seq data should be deposited in a public database.

52 **Reply:** Page 22. We have deposited the raw sequencing data and associated meta data on the
53 European Nucleotide Archive under the study accession number: PRJEB38464
54 ([url{http://www.ebi.ac.uk/ena/data/view/PRJEB38464}](http://www.ebi.ac.uk/ena/data/view/PRJEB38464)).
55

56
57 10. Some of the data here might reveal additional genetic differences if the read pairs from the
58 paired-end sequencing were analyzed together. For example, clusters of read pairs with longer than
59 expected gaps in the reference genome would suggest a deletion in the sequenced individual.
60

1
2
3 Insertions in the sequenced individual might lead to some read pairs with only one read being
4 mapped (because the other read is in the inserted sequence). Such analyses should be performed to
5 identify genome differences, which can be considered for population structure.

6 Reply: Page 8. This could be an interesting approach but we would be concerned that with only 100
7 bp “loci” for RAD data, reliably interpreting insertions and deletions would be quite challenging. A
8 limitation of RAD sequencing is that allelic dropout is impossible to distinguish completely from lack
9 of coverage. Our filtering decisions were aimed at reducing these types of errors, as described in
10 O’Leary et al. (2018). We have added to our methods to clarify the filtering decisions that were also
11 queried by reviewer 3. Lines 195--202.

12
13
14 11. Key statistics of the RAD-seq data should be reported in a supplemental table for each individual:
15 (1) number of raw read pairs; (2) number of read pairs with mapping results with expected gap size;
16 (3) number read mapped to the genome; (4) SNPs before filtering; (5) number of reads with SNPs;
17 (6) number of SNPs after filtering; (7) range of sizes of RAD loci; (8) number of filter SNPs on each
18 chromosomal arms (18 numbers), and possibly others.

19
20 Reply: We have included details about the number of reads, the number of reads that mapped to
21 the genome, the number of SNPs and missing SNPs in supplementary tables separately. In addition,
22 we have provided information in the supplementary from the output of running stacks to enable the
23 reader to find any additional information that they may wish to see.

24
25
26 12. Some SNPs might be false positives because they are actually from two or more copies of nearly
27 identical sequences in the genome. Sometimes sequenced individual might have more copies than
28 the reference genome. Therefore, de novo assembly (above) can also help to cover such mistakes in
29 SNP identification.

30 Reply: The filtering used was designed to reduce inclusion of duplicate loci. The minimum stack
31 depth used was 5 (default is 2), which is in the middle range of depths usually considered for
32 excluding loci (Paris *et al.* 2017 *Methods in Ecology and Evolution*: 8 1360-1373). This was used to
33 minimise the number of paralogues (a higher value could lead to filtering out informative loci). Using
34 a maximum observed heterozygosity of 0.7 also helps to remove potential paralogues as no more
35 than 70% of the individuals (dataset 1 and 2) or populations (dataset 3) can be heterozygous for
36 each locus.

37
38
39 13. *B. oleracea* and other related species have shared a genome triplication, and many paralogs are
40 still retained. There are other duplications. These homologous sequences might also differ in copy
41 number among individuals, and be additional challenges to correct identification of SNPs.
42 For example, if *B. oleracea* has two highly similar genes (A and B), but one of them is lost in the
43 reference genome or not sequenced (B), and the allele of the second copy (B) in some wild
44 individual could be treated as allelic to A if the true A allele is not detected by the RAD-seq.
45 Mapping of both of the paired reads can reduce such mistakes, and careful analysis of multiple reads
46 of the same locus can also help.

47
48 Reply: Page 16. This would be very useful to consider if we were attempting to reconstruct the
49 genome sequences but ddRAD data are not the most appropriate for resolving paralogs, given the
50 filtering to include only a single SNP per locus, a standard approach for population genetic analyses
51 (Pritchard et al. 2000), and the short length of the “loci” between restriction sites (100 bp). Instead,
52 duplicate loci are excluded by having a minimum stack depth of 5, a maximum observed
53 heterozygosity of 0.7 and a minimum minor allele frequency of 0.01. In addition, we only consider
54 those loci that are present in at least 60% of the individuals in datasets 1 & 2, and in 50% of the
55 populations in dataset 3.

56
57
58 We have added a paragraph to the discussion that more explicitly describes the limitations of ddRAD
59 sequencing for enabling comparisons with published genome sequences and other types of
60

1
2
3 genotype by sequencing data (lines 411--427). This would be particularly problematic for complex
4 Brassicaceae genomes, precisely due to the triplication. However, for population genetic analyses,
5 the filtering applied here was conservative enough that these issues should be reduced. Although
6 the excess heterozygosity observed could be explained by introgression, we also more explicitly
7 acknowledge that some could be due to merging of duplicate loci. However, the consistency in this
8 excess across populations suggests that interpretations about relative diversity and population
9 structure should still be meaningful.
10

11
12 14. Different methods for phylogenetic analysis should be use to demonstrate the strength of the
13 results. Methods that consider conflicts among the loci are especially valuable.

14 Reply: Since we are considering population genetic variation, we are not expecting bifurcating trees
15 appropriate to a rigorous phylogenetic analysis. Instead, RAXML was used for clustering, to visualise
16 variation within and among populations. We also used Splitstree for visualisation, which explicitly
17 considers conflicts among loci by presenting alternative pathways.
18

19
20 15. In addition, different subsets of the data can be used to generate the phylogeny.

21 Reply: We hadn't meant to imply that a phylogenetic approach was appropriate for the within-
22 species population-level data presented. We checked the manuscript to make sure that we hadn't
23 referred to a phylogeny. The cluster analysis presented is for visualisation of relative variation within
24 and between populations. The level of admixture apparent in figure 1-iii (STRUCTURE plot) clearly
25 demonstrates why a phylogenetic approach would not be appropriate.
26
27

28
29 Reviewer: 3
30

31 Comments to the Author

32 The study uses ddRAD genotyping in the domesticated plant species Brassica oleraceae, collected at
33 24 localities along the coasts of the UK and Spain, to characterize the population structure of wild
34 'populations' and to explore environment-genotype associations across varying climatic conditions.
35 In general, the manuscript reads well. The data are newly generated for this study. However, there
36 are some parts of the methods and results that need clarification (outlined below) as I could not
37 understand important aspects of the research. Although I think the results are relevant, I don't think
38 the manuscript is appropriate for the broad readership of Molecular Ecology.
39
40

41 Major comments

42 -I have major reservations about the sample sizes used in each 'population' as well as the
43 designation of localities as populations. In the manuscript, each 'population' has 4 samples, but
44 population-level statistics are very sensitive to sampling error at this small sample size, and I worry
45 that they are unreliable. The authors should demonstrate that their analyses are robust to their
46 sample sizes to make their conclusions more convincing. A couple of possibilities would be to drop
47 the population sample size from 2 to 4 to see if overall patterns hold up, or using simulated data to
48 prove that the methods are reliable with a sample size of 4. Besides, the authors treat localities as
49 populations. A typical biological definition of a population is a group of interbreeding individuals that
50 share time and space (Hedrick PW (2000) Genetics of Populations, 2nd edn. Jones and Bartlett,
51 Sudbury, Massachusetts.), and most definitions involve some type of interbreeding. One of the goals
52 of population genetics is to identify what is a population and how many are present. Many statistical
53 tests were developed to identify populations. We need to be precise with our terminology. It would
54 be better if you used sampling location, or site, in place of population.
55

56 Reply: Page 6. We agree with the reviewer that the sample sizes seem small initially. However, a
57 recent study into the minimum sample sizes required for population genomic analyses (Nazareno et
58 al 2017; Molecular Ecology Resources 17(6): 1136--1147) found that very small sample sizes are
59
60

1
2
3 useful for various population statistics when large numbers of SNPs are available. For example, they
4 found that with ≥ 1500 SNPs from two individuals, F_{ST} and H_o could be accurately estimated, and
5 that with four individuals H_e could also be reasonably estimated. Therefore, although larger sample
6 sizes would be more ideal for estimating H_e , the sample sizes used here (2-4 individuals per locality
7 and $>13,000$ SNPs) provide enough information to obtain an overview of the genetic structuring and
8 shared ancestry. We have also previously used this approach to investigate population structure in
9 wild Brassicaceae (genus *Arabidopsis*). We have added these additional supporting references to the
10 methods (lines 153--156).
11
12

13 -There is almost no information given about the ddRAD loci and SNPs used for the analyses. For
14 instance, what is the distribution of allele frequencies of the SNPs? Without this information, it is
15 impossible to gauge whether the analyses are appropriate for the data. Furthermore, no summary
16 statistics were provided for these data. Please provide tests for Hardy-Weinberg equilibrium, linkage
17 disequilibrium, and tests for allelic dropout. Allelic dropout is a serious issue for RADseq data.
18 Further, you state the SNPs are unlinked but do not say you tested for this. Please present the tests,
19 such as LD, that told you these are unlinked.
20

21 Reply: Page 8. The use of "unlinked" was not descriptive enough. We meant to use it in the context
22 of being unlinked within a locus, by only retaining one SNP within each locus for dataset 2, which is
23 used in the population structure analysis. We have made the filtering steps clearer in our methods
24 section (lines 195--202), which should remove some of the confusion.
25

26 We originally checked for HWE with PLINK but failed to acknowledge this. We have re-run the
27 population-level analyses for dataset 3 to include the estimation of HWE within stacks. Table 3 has
28 been updated based on this re-run, but the results did not drastically change and therefore did not
29 alter our conclusions.
30

31 Minor comments

32 -The authors need to clarify how many populations and individuals were used in this study (e.g., 20
33 or 24 populations? 96 or 144 individuals?).
34

35 Reply: Page 6, line 162. The 144 was what was originally attempted to be sequenced, but
36 unfortunately sequencing was not successful for many of these samples/populations. We have
37 changed this to just include those populations where some individuals were successfully sequenced:
38 96 individuals (4*24). Furthermore, we have added supplementary information about the number of
39 reads from these 96 individuals, including those excluded from downstream analyses.
40

41 -More details about the ddRADseq genomic library preparation should be given.
42

43 Reply: Pages 6 and 7. Additional information has been included in the methods lines 165--178.
44 Specifically: "Double-digest RADseq libraries were made using a modification of the method in Wu et
45 al. 2016 that allowed NexteraXT indexes (Illumine Corp., USA) to be used for multiplexing samples. In
46 addition, an RYRY spacer was inserted in the adapter 3' of the Illumina sequencing primer annealing
47 site to provide additional complexity at the start of read 1 immediately before the Sac1 sticky end.
48 For each sample 400 ng DNA was fully digested with Sac1 and Mse1 restriction endonucleases and
49 purified using Ampure XP beads. Illumina compatible i5 adapters were designed to ligate to the at
50 the AGCT-3' sticky end left after Sac1 digest, and Illumina compatible i7 adapters were designed to
51 ligate to the 5'-TA overhangs remaining after Mse1 digest. Adapter-ligation excess adapters were
52 removed using Ampure XP beads. DNA fragments were amplified by 12 cycles of indexing PCR,
53 purified, size selected (inserts 330-670 bp) and validated using a TapeStation D1000 HS Screentape
54 (Agilent Technologies Ltd). Libraries were equimolar pooled and the pool concentration was
55 calculated after qPCR. Libraries were denatured, diluted and sequenced with 125bp paired-end
56 reads on Illumina HiSeq 2500 using SBS High Output reagents v4 (Illumina Corp., USA)."
57
58
59
60

1
2
3
4 -The authors also need to provide more details about the parameters used in populations and better
5 explain why they were used to obtain each data set. Also, authors need to justify why they choose
6 so low values for some parameters in populations (e.g., -r and -p values). How these values can
7 affect your results?
8

9 **Reply: Page 8. The parameters used in populations were chosen to balance the amount of missing**
10 **data against the number of SNPs retained keeping in mind the need to reduce inclusion of duplicate**
11 **loci. We have added a paragraph to the methods to explain this more fully, lines 195-202:**

12 **“This filtering was designed to reduce the inclusion of duplicate loci and balance the amount of**
13 **missing data with the number of informative loci (Andrews et al 2016). A minimum stack depth of**
14 **five is higher than the default of two, but within the recommended range (Paris et al 2017), and**
15 **helps to remove potential paralogues. Spurious SNPs are avoided by using a minor allele frequency**
16 **of >0.01 (Marandel et al 2020), and the combination of a maximum observed heterozygosity of 0.7**
17 **(70 % of the individuals or populations can be heterozygous for each locus) which are present in**
18 **either 60 % of individuals (datasets 1 and 2) or 50 % of the populations (dataset 3) retains loci that**
19 **have been successfully genotyped across individuals, but are not completely heterozygous.”**
20
21

22 -I didn't understand why authors used cutadapt to demultiplex and to trimm reads since the
23 process_radtags pipeline in Stacks does these procedures.

24 **Reply: We used cutadapt instead of process_radtags for demultiplexing and trimming because the**
25 **library preparation included an additional spacer in the adapter 3' Illumina primer. This was to allow**
26 **for additional complexity in the libraries and improve the cluster identification and registration on**
27 **the HiSeq. It was more straightforward to remove these using cutadapt at the time of analyses.**
28
29

30 -I would like to know if authors had some issues with the ddRADseq genomic library, since so few
31 reads (i.e., 16,894,310) were obtained for the entire sampling size (n=76 individuals). It is correct? If
32 yes, the average reads per individual is 222,293 and not 1,534,680.

33 **Reply: Page 10. Thank you for point this out, it was a typo which meant that the total number was an**
34 **order of magnitude out. This has been corrected (line 249).**
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60