



Deposited via The University of Sheffield.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/161665/>

Version: Accepted Version

Article:

Kim, J.-M., Santure, A.W., Barton, H.J. et al. (2018) A high-density SNP chip for genotyping great tit (*Parus major*) populations and its application to studying the genetic architecture of exploration behaviour. *Molecular Ecology Resources*, 18 (4). pp. 877-891. ISSN: 1755-098X

<https://doi.org/10.1111/1755-0998.12778>

This is the peer reviewed version of the following article: Kim, J-M, Santure, AW, Barton, HJ, et al. A high-density SNP chip for genotyping great tit (*Parus major*) populations and its application to studying the genetic architecture of exploration behaviour. *Mol Ecol Resour.* 2018; 18: 877– 891, which has been published in final form at <https://doi.org/10.1111/1755-0998.12778>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions.

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

1 A high density SNP chip for genotyping great tit (*Parus major*) populations and
2 its application to studying the genetic architecture of exploration behaviour

3

4 J-M Kim^{1,6}, AW Santure^{1,7}, HJ Barton¹, JL Quinn², EF Cole³, Great Tit HapMap Consortium, ME
5 Visser⁴, BC Sheldon³, MAM Groenen⁵, K van Oers⁴ & J Slate¹

6

7 Addresses

8 1. Department of Animal & Plant Sciences, University of Sheffield, Sheffield, S10 2TN, UK

9 2. School of Biological, Earth and Environmental Science (BEES), University College Cork, Distillery
10 Fields, North Mall, Cork, Ireland

11 3. Edward Grey Institute, Department of Zoology, University of Oxford, Oxford, OX1 3PS, UK

12 4. Department of Animal Ecology, Netherlands Institute of Ecology (NIOO-KNAW), Wageningen,
13 Netherlands

14 5. Wageningen University and Research - Animal Breeding and Genomics, Netherlands

15 6. Department of Animal Science and Technology, Chung-Ang University, Anseong, Gyeonggi-do, 456-
16 756, Republic of Korea.

17 7. School of Biological Sciences, University of Auckland, Private Bag 92019, Auckland 1142, New
18 Zealand

19

20 Keywords: GWAS, Axiom, exploration behaviour, personality, CNV

21

22 Corresponding author: Jon Slate (j.slate@sheffield.ac.uk)

23

24 Abstract

25 High density SNP microarrays ('SNP chips') are a rapid, accurate and efficient method for
26 genotyping several hundred thousand polymorphisms in large numbers of individuals. While
27 SNP chips are routinely used in human genetics and in animal and plant breeding, they are
28 less widely used in evolutionary and ecological research. In this paper we describe the
29 development and application of a high density Affymetrix Axiom chip with around 500 000
30 SNPs, designed to perform genomics studies of great tit (*Parus major*) populations. We
31 demonstrate that the per-SNP genotype error rate is well below 1% and that the chip can
32 also be used to identify structural or copy number variation (CNVs). The chip is used to
33 explore the genetic architecture of exploration behaviour (EB), a personality trait that has
34 been widely studied in great tits and other species. No SNPs reached genome-wide
35 significance, including at *DRD4*, a candidate gene. However, EB is heritable and appears to
36 have a polygenic architecture. Researchers developing similar SNP chips may note: (i) SNPs
37 previously typed on alternative platforms are more likely to be converted to working assays,
38 (ii) detecting SNPs by more than one pipeline, and in independent datasets, ensures a high
39 proportion of working assays, (iii) allele frequency ascertainment bias is minimised by
40 performing SNP discovery in individuals from multiple populations and (iv) samples with the
41 lowest call rates tend to also have the greatest genotyping error rates.

42

43 Introduction

44 It is now becoming commonplace to sequence and assemble the genomes of organisms that
45 have been the focus of ecological research but are not classical genetic model organisms
46 (Brawand *et al.* 2014; Colbourne *et al.* 2011; Ellegren *et al.* 2012; Hu *et al.* 2011; Jones *et al.*
47 2012; Lamichhaney *et al.* 2015; Soria-Carrasco *et al.* 2014). While assembled genomes are
48 undoubtedly essential tools for understanding topics in evolutionary and ecological genetics,
49 in taxa with moderate to large genomes the cost of sequencing the full genomes of hundreds
50 or thousands of individuals remains prohibitive for the majority of laboratories, and beyond
51 the budget of even very large grants. Thus, analytical techniques that require large sample
52 sizes, such as quantitative trait locus (QTL) linkage mapping / genome-wide association
53 studies (GWAS) (Visscher *et al.* 2017), molecular quantitative genetics (Gienapp *et al.* 2017a;
54 Jensen *et al.* 2014) and studies that utilise realised relatedness / inbreeding coefficients
55 (Powell *et al.* 2010) are reliant on alternative technologies. Broadly, these can be categorised
56 into two approaches; (i) genotyping-by-sequencing (GBS) methods (Davey *et al.* 2011) such as
57 restriction-site associated sequencing (RAD-seq) (Hohenlohe *et al.* 2010) and double-digest
58 RAD-seq (ddRAD-seq) (Peterson *et al.* 2012) and (ii) SNP microarray ('SNP chip') methods
59 (Spencer *et al.* 2009; Syvanen 2001), where a set of known SNPs are probed on chips
60 manufactured by providers such as Illumina (Shen *et al.* 2005) and Affymetrix (Matsuzaki *et*
61 *al.* 2004).

62

63 GBS-approaches while perhaps cheaper, are more technically demanding, both in terms of
64 laboratory work, and in post-sequencing processing of NGS data (Bajgain *et al.* 2016; Miller *et*
65 *al.* 2012; Robledo *et al.* 2017). Furthermore, the sites that are typed are typically not known
66 in advance, and call rates can vary widely between different SNPs. SNP chips are more
67 expensive, but tend to have higher call rates per SNP, and specific target SNPs can be
68 included in chip design. In addition, the same SNPs are typed in every individual, which is not
69 the case for GBS approaches (Bajgain *et al.* 2016). A disadvantage of SNP chips is
70 ascertainment bias (Bajgain *et al.* 2016; Miller *et al.* 2012). Because SNPs have to be
71 discovered before they are designed to be on a chip, there is usually bias towards the
72 inclusion of SNPs with higher minor allele frequencies (MAF) on the chip. For some types of

73 analyses (e.g. GWAS) this is not necessarily a disadvantage, because statistical power is
74 greater for SNPs with higher MAF. However, ascertainment bias is clearly a problem for tests
75 that require an accurate description of the site frequency spectrum in different genomic
76 regions (Albrechtsen *et al.* 2010) e.g. tests that aim to detect signatures of selection such as
77 Tajima's D. Thus, the optimal method for genotyping many individuals can depend on the
78 question being addressed, the laboratory and bioinformatics experience of the user and the
79 laboratory budget.

80

81 The great tit (*Parus major*) is a model vertebrate system in evolutionary ecology because this
82 passerine bird readily breeds in nest boxes (making it possible to identify parents and
83 offspring and thus build pedigrees), it has a short generation time and large broods, and it is
84 widely distributed across Europe, Western Asia and parts of the Middle East (Perrins 1979).
85 Longitudinal studies (Kluyver 1951; Lack 1964) of great tits have informed researchers about
86 classic topics in evolutionary and behavioural ecology (Lack 1968) including mating systems
87 and reproductive decisions (Smith *et al.* 1989), the frequency (Harvey *et al.* 1979) and
88 importance of dispersal (Garant *et al.* 2005; Postma & van Noordwijk 2005), adaptation to
89 climate change (Charmantier *et al.* 2008; Nussey *et al.* 2005; Visser *et al.* 1998), the study of
90 personality traits (Dingemanse *et al.* 2004; Groothuis & Carere 2005; Van Oers & Naguib
91 2013), innovativeness and cognition (Cole *et al.* 2012; Quinn *et al.* 2016; Titulaer *et al.* 2012),
92 social learning (Aplin *et al.* 2015; Aplin *et al.* 2012), and understanding how quantitative
93 genetic variation is maintained in natural populations (McCleery *et al.* 2004). In more recent
94 years, great tits have become the focus of molecular genetic studies exploring the genetic
95 architecture of quantitative traits (Gienapp *et al.* 2017b; Robinson *et al.* 2013; Santure *et al.*
96 2013; Santure *et al.* 2015), phylogeography (Kvist *et al.* 2003; Lemoine *et al.* 2016), fine-scale
97 genetic structure and dispersal (Garroway *et al.* 2013; Radersma *et al.* 2017), the efficacy,
98 nature and relative occurrence of positive and purifying selection (Corcoran *et al.* 2017;
99 Gossmann *et al.* 2014) and immunogenetics (Sepil *et al.* 2013; Sepil *et al.* 2012). Much of this
100 work has been facilitated by a SNP chip containing probes for around 10,000 SNPs, of which
101 around 6,000 are polymorphic and reliably scoreable (Van Bers *et al.* 2012). This '10K chip'
102 has been used in QTL and GWAS mapping studies and to construct a great tit linkage map
103 (van Oers *et al.* 2014) which led to insights into the nature of sex-differences in

104 recombination rate (heterochiasmy). The linkage map was in turn used to help assemble the
105 great tit genome (Laine *et al.* 2016).

106 While the 10K SNP chip has helped provide insight into the architecture of some quantitative
107 traits, it also suffers from some important limitations (Santure *et al.* 2015). The most
108 important of these is that the marker density (~ 1 SNP per 20Kbp) is too low for most of the
109 genome to be adequately 'tagged' by typed SNPs that are in strong linkage disequilibrium
110 (LD) with untyped sites. Furthermore, molecular quantitative genetic approaches such as
111 chromosome partitioning (Yang *et al.* 2011) or regional heritability mapping (Nagamine *et al.*
112 2012), where markers are used to measure between-individual relatedness in specific
113 genomic regions, typically require a much higher marker density than is afforded by the 10K
114 chip (Berenos *et al.* 2014).

115 To overcome the low power of the 10K chip, and to provide better resolution in association
116 studies, outlier detection tests and molecular quantitative genetic analyses we have
117 developed a high density (HD) chip with probes for over 600 000 SNPs. In this paper we
118 describe the development of this great tit HD SNP chip. The chip can also be used to detect
119 the presence of structural variation or copy number variants (McCarroll & Altshuler 2007) in
120 the great tit genome. We demonstrate an application of the HD chip, using a behavioural
121 trait, to showcase how the genetic architecture of phenotypic variation can be estimated. It is
122 hoped that the methods and lessons described in this paper will serve as a useful guide to
123 researchers developing high density SNP chips in other organisms.

124 **Methods**

125

126 **DNA sequencing**

127 To identify SNPs to include on the chip, whole genome resequencing was performed on 30
128 birds. Ten of the birds were from the long term study population at Wytham Woods, Oxford,
129 UK (51°46' N, 1°20' W), and the remaining 20 were from locations across a wide area of
130 Europe (Fig. S1), collected as part of the Great Tit HapMap Project. The sequencing is
131 described elsewhere (Laine *et al.* 2016), but briefly, samples were sequenced on an Illumina
132 HiSeq 2000 platform at The Genome Institute, Washington University. Sequencing was
133 paired-end, with insert sizes 300 bp and a read length of 100 bp. Each bird was sequenced to
134 ~10x coverage. Note that one of the samples used in this paper, from near to Zurich in
135 Switzerland (population #27 in Fig. S1), was not used in the genome assembly paper (Laine *et*
136 *al.* 2016), because coverage was lower than for other samples (~5x). The Zurich sample is
137 included in the NCBI sequence read archive submission (SRP066678).

138

139 **SNP Discovery**

140 SNP discovery was performed in several steps, with the aim of identifying markers that are
141 polymorphic across multiple great tit populations, with minimal ascertainment bias towards
142 populations where the SNPs were initially discovered. Paired-end reads were filtered and
143 trimmed with the FASTX-Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/) using a length of
144 80 bp and quality score of 20 as minimum cut-off scores to remove low-quality reads. The
145 remaining reads from each individual were mapped onto the great tit reference genome
146 v1.03 with the MEM algorithm of the Burrows-Wheeler Aligner (Li & Durbin 2009). The
147 aligned sequence reads on the genome were stored as individual BAM files. Using VCFtools
148 (Danecek *et al.* 2011), the BAM files were filtered to a minimum quality score of 20 and read
149 depth of 5.

150 Following alignment of reads to the great tit genome, a combination of different SNP
151 discovery algorithms and different strata of the dataset were used, summarised in Figure 1.

152 SNPs were independently called using the ANGSD v0.549 (Korneliussen *et al.* 2014),
153 SAMTools v0.1.19 (Li 2011; Li *et al.* 2009) and GATK v2.4 (DePristo *et al.* 2011; McKenna *et al.*
154 2010) packages. Parameter settings are reported in Table S1. SNPs were called either from
155 the 10 UK birds, the 20 mainland European birds, or the combined dataset of 30 birds. SNPs
156 called from the different software/datasets were then compared (Figure 1) and a set of
157 ~1.4M SNPs that were common to all SNP discovery softwares and all datasets were
158 considered for inclusion on the SNP chip. VCFtools was used to filter out SNPs with minor
159 allele frequency (MAF) less than 0.05 and call quality less than 50. SNPs that were predicted
160 to be within 30 bp of each other were filtered out because it was likely that the presence of
161 one SNP would adversely affect the ability to successfully genotype the other(s), due to
162 inefficient or biased hybridisation of allele-specific oligonucleotides. SNPs prone to this form
163 of possible typing error are known as Off Target Variants (OTV) in the Affymetrix genotype
164 calling workflow (see below). A total of 1,213,160 SNPs passed all of these filtering criteria
165 (Figure 1).

166

167 SNP selection

168 The SNP discovery phase of the work identified more SNPs than could be included on the
169 chip. To prioritise which SNPs to use on the chip, the following criteria were used:

170 1) 'Top priority' SNPs were those that had been successfully typed on the lower density
171 10K chip described in earlier work (Van Bers *et al.* 2012) or had been discovered in
172 the SNP discovery pipeline described above *and* were discovered during the
173 construction of the earlier 10K chip but not included on it (Santure *et al.* 2011; van
174 Bers *et al.* 2010). 6,773 SNPs that were typed on the original chip and a further 9,713
175 SNPs that were discovered but not included on the 10K chip were included in the 'Top
176 priority' set.

177 2) A list of candidate genes were identified that could potentially explain variation in
178 ecologically relevant traits such as personality traits (Fidler *et al.* 2007; van Oers *et al.*
179 2004) and timing of breeding (Visser *et al.* 2003). A list of candidate genes and
180 putatively associated traits is provided in Table S2. At the time the chip was being
181 designed, the great tit genome was not annotated. Therefore, to identify the location

182 of the candidate genes on the great tit genome, the cDNA sequence of the candidate
183 gene in zebra finch (*Taeniopygia guttata*), another passerine, chicken (*Gallus gallus*),
184 or if none of those were available, human or mouse, was downloaded from NCBI and
185 the location on the great tit genome was identified by BLAST search. The start and
186 end point of the gene was identified and SNPs were considered for inclusion if they
187 were within any part of the gene. 654 (of which 28 were also 'Top Priority' SNPs) from
188 110 genes were chosen for inclusion on the chip.

189 3) The remaining SNPs were selected based on how likely they were to be convertible to
190 a working and scoreable assay on the chip. The list of SNPs and their flanking
191 sequences were sent to the Affymetrix bioinformatics team who used their *in silico*
192 design tool to model the probability (termed the 'P convert design score') of the SNP
193 converting to a working assay. The software uses the SNP bases and its flanking
194 sequence, and considers factors such as GC content and the predicted amount of
195 non-specific hybridisation to other (non-target) genomic regions. Following this
196 process, SNPs with a P convert design score >0.69 were retained for inclusion on the
197 chip. This threshold compares favourably to those used in the design of HD chips for
198 chicken (Kranis *et al.* 2013), catfish (Liu *et al.* 2014), and water buffalo (Iamartino *et*
199 *al.* 2017), where thresholds of 0.20, 0.50 and 0.60 were used respectively.

200 An Axiom myDesign high density chip was manufactured by Affymetrix. A total of 610 970
201 SNPs were included on the final design, of which 17 122 were from criteria 1 or 2 and the
202 remainder were from criteria 3. The genomic distribution of attempted SNPs are described in
203 Table S3 and Fig. S1.

204

205 Genotyping

206 Genotyping was performed on a Gene Titan platform at Edinburgh Genomics. A total of 21
207 plates, each with up to 96 samples, were typed (2016 available slots). Across the 21 plates, 9
208 negative controls were included. All plates contained at least one duplicate sample to aid
209 with estimation of error rate. 1073 typed samples were from the Wytham Woods population.
210 The remainder of the total 2007 birds came from a number of study sites (Table 1, Fig. S1)
211 from across the species range in Europe and Asia, and were provided by members of the

212 Great Tit HapMap Consortium, either as pre-extracted DNA, or more usually as blood
213 samples in Queen's storage buffer or ethanol. DNA was extracted using an ammonium
214 acetate precipitation method (Bruford *et al.* 1998) and DNA quality and quantity measured
215 using picogreen on a fluorometer. 1,696 samples were at a concentration exceeding 50ng/ul,
216 while 89 were at concentrations lower than 20ng/ul. All except 33 samples passed the
217 manufacturer's recommendation of 200ng of DNA. 13 Japanese tit (*Parus minor*) birds were
218 genotyped, as well as 9 putative *P. major* / *P. minor* hybrids. Abel, the male used as the
219 reference bird for the great tit genome assembly (Laine *et al.* 2016), was typed four times
220 (two replicates on two different plates). SNP genotype calling was performed using the
221 Ps_Metrics and Ps_Classification functions within the Affymetrix Axiom Analysis Suite
222 1.1.0.616. Samples with dish QC < 0.82 or call rates <0.95 were discarded, as were SNPs with
223 call rates <0.97 or those identified as containing Off-Target Variants (OTVs).

224 **Quantifying Genotyping Error Rate**

225 Genotyping errors were estimated in two ways. First, the replicated samples meant that the
226 proportion of inconsistent genotypes between different typing attempts of the same bird
227 could be estimated. The error rate was obtained from the Z2 score - the proportion of SNPs
228 at which two individuals (replicates) share both alleles identically-by-descent - reported by
229 the --genome command in Plink 1.9 (Chang *et al.* 2015). Second, genotypes from the SNP
230 chip were compared with the whole genome resequencing SNP calls for 28 birds that were
231 successfully genotyped and sequenced to ~10x coverage (Laine *et al.* 2016). Note that
232 discrepancies between chip and resequencing SNP genotypes can arise either because the
233 SNP chip genotype is wrong, or because the SNP call from the resequencing is wrong.
234 Therefore, comparison between the resequencing and the SNP chip genotyping provides an
235 upper limit on the genotyping error on the SNP chip. Concordance between the chip and the
236 resequencing data was determined using the GenotypeConcordance tool implemented
237 within GATK, after SNPs with Genotype Quality Scores <30 were filtered from the
238 resequencing dataset.

239

240 **Copy number variant (CNV) detection**

241 CNVs were detected using the PennCNV software (Wang *et al.* 2007). PennCNV input files of
242 the 996 birds from the Wytham Woods population were prepared using the Axiom Analysis
243 Suite's CNVTool and probe intensities from all SNPs. PennCNV uses two parameters from the
244 SNP genotyping, the logR ratio and the B allele frequency, to identify genomic segments
245 containing SNPs indicative of copy number variation. The logR ratio is a measure of signal
246 intensity. SNP assays in individuals with extra copies of a genomic region (duplications)
247 should generate higher intensity signals, while SNPs in individuals with fewer than two copies
248 of a genomic segment (deletions) should generate lower intensity signals. The B allele
249 frequency measures the relative signal intensity of the two possible alleles at each SNP.
250 Ratios that are inconsistent with allele call ratios of 2:0 (i.e. A allele homozygote), 1:1 (i.e.
251 heterozygote) or 0:2 (i.e. B allele homozygote) are indicative of departures from two copies
252 of that nucleotide (i.e. the normal diploid state) being present in the sample. For example, an
253 individual with a duplication at a CNV site on one chromosome, would have three copies in
254 total, meaning the ratios of alleles A:B could be 1:2 or 2:1, which is impossible when two
255 copies are present. CNVs called by PennCNV were retained and converted to Plink format
256 using the perl script `penncnv_to_plink`
257 (www.openbioinformatics.org/penncnv/download/penncnv_to_plink.pl). The plink
258 commands `--cfile` `--cnv-overlap` and `--cnv-seglist` were used to generate a list of all CNVs,
259 identify overlapping CNVs, estimate CNV frequencies and summarise the CNVs present in
260 each individual (.cnv.indiv file).

261 Additional CNV analyses included (i) an examination of two replicates of the reference
262 genome bird, Abel, and (ii) CNV calling using nine father-mother-offspring trios from the
263 Wytham Woods population. As with the analysis of all Wytham Woods birds, the PennCNV
264 command `detect_cnv_pl` was used, only with the `-trio` argument included. In principle,
265 detected CNVs are more likely to be reliable calls if they are observed to be inherited in a
266 Mendelian fashion.

267

268 Genetic architecture of a personality trait

269 The chip was used to explore the genetic architecture of Exploration Behaviour in a novel
270 environment (EB), a personality trait linked to aggression, risk-taking and dispersal in great

271 tits (Quinn *et al.* 2009). EB is known to be heritable (Dingemanse *et al.* 2004; Drent *et al.*
272 2003; Quinn *et al.* 2009; Santure *et al.* 2015) and it has also been the focus of candidate gene
273 studies, especially at the Dopamine D4 receptor (*DRD4*) gene (Fidler *et al.* 2007; Korsten *et al.*
274 2010), following the first report that *DRD4* could affect novelty-seeking behaviour in humans
275 (Ebstein *et al.* 1996). The protocol for measuring EB is described in detail elsewhere (Cole &
276 Quinn 2014; Quinn *et al.* 2009). Briefly, wild birds were captured during February-March
277 (2005) or September-March (2006-2009) and assayed in a novel environment room at
278 Wytham Woods field station. For the purposes of the downstream genetic analyses we used
279 the same measure of EB as that used in previous studies. Briefly, the first principal
280 component (PC1) of 12 behavioural measures was treated as the EB score. PC1 was square-
281 root transformed prior to genetic analysis and a single value for each individual was obtained
282 by fitting a linear mixed model with the terms ID, year, days after September 1st, and assay
283 number of that individual all included as predictors. Details are described elsewhere (Quinn
284 *et al.* 2009). Several aspects of EB genetics were explored. First, we performed a genome-
285 wide association study (GWAS) using the Grammar method (Aulchenko *et al.* 2007a),
286 implemented in GenABEL (Aulchenko *et al.* 2007b). Grammar accounts for the possibility of
287 test statistic inflation caused by relatives in the dataset by fitting a realised genome-wide
288 relationship matrix estimated from the SNP data as a random effect. The residual from the
289 random model was used as the phenotype. In addition, genomic correction was performed
290 by estimating lambda, the slope of observed chi square values on expected chi square values,
291 and dividing all tests statistics by lambda before estimating nominal P-value. Genome-wide
292 statistical significance was estimated by permutation test, using the GenABEL mmscore
293 command and 1000 permutations of the data. The GWAS was performed on a total of 415
294 birds from Wytham Woods. All Z-linked SNPs and any autosomal SNPs with MAF <0.05 or
295 significant departures from Hardy-Weinberg Equilibrium ($P < 1 \times 10^{-5}$) were filtered from the
296 dataset leaving a total of 459 502 autosomal SNPs.

297 In addition to the GWAS, an additional analysis of the same dataset fitted all SNPs
298 simultaneously, in one model. Here, the objective was to estimate the proportion of
299 phenotypic variation explained by each SNP, in order to understand aspects of the trait
300 architecture such as the heritability, the number of SNPs in linkage disequilibrium with causal
301 variants and the distribution of effect sizes of those SNPs. The BayesR method (Erbe *et al.*

302 2012), whereby it is assumed that the SNPs causing phenotypic variance are drawn from a
303 mixture of different effect size distributions, was used to model the genetic architecture of
304 EB. The BayesR package (Moser *et al.* 2015) was used to run the analyses, with default
305 settings of 4 distributions, with mean effect sizes of 0.01, 0.001, 0.0001 or 0 of the
306 phenotypic variation. The program was run for 50 000 iterations of an MCMC chain, with the
307 first 20 000 iterations treated as burn-in, and every 10th chain after that being sampled,
308 giving a total of 3000 samples of the chain. Priors for V_A and V_E were specified using an
309 inverted chi-squared distribution with scale parameters of 0.033 and 0.117 respectively, each
310 with 4 degrees of freedom. These values give a prior heritability of around 0.20 which is
311 consistent with pedigree-based estimates of EB in the Wytham Woods population (Quinn *et*
312 *al.* 2009; Santure *et al.* 2015). Note that setting the priors so that V_A and V_E were identical
313 (i.e. the heritability was 0.5) gave almost identical posterior estimates, so the genetic
314 architecture does not appear to be sensitive to the priors.

315

316 Results

317 Summary Statistics

318 Following genotype calling and quality control steps, a total of 1 846 samples typed at 502
319 685 SNPs were retained for analysis. A summary of the different types of SNP category is
320 provided in Table 2. Samples that contained less than the recommended 200ng of DNA were
321 more likely to fail than those with >200ng of DNA; 9/33 failures versus 140/1962 failures
322 (Fisher's Exact Test: Odds ratio = 4.87, 95% CI 1.95-11.12, P = 0.0005). However, among
323 samples that passed quality control, there was no relationship between the call rate and the
324 amount of DNA present in the sample ($F_{1,1844} = 0.942$, P = 0.33). SNPs that had been
325 previously typed on the 10K chip were more likely to be converted to a successfully typed
326 SNP, and to pass QC checks. For previously typed SNPs the conversion rate was 5924/6773
327 (0.87) compared to 496 826 / 604 197 (0.82) for unvalidated SNPs; Fisher's Exact Test odds
328 ratio 1.51, 95% CI = 1.40-1.62, P = 0.0006. However, SNPs that were discovered during both
329 the construction of the 10K chip and of the HD chip but were not typed on the 10K chip
330 actually had a lower conversion success rate, 7807/9713 (0.80), than SNPs that were only
331 discovered during HD chip construction, 489 019 / 594 484 (0.82); Fisher's Exact Test: Odds
332 ratio = 0.88, 95% CI = 0.84-0.93, P = 2.0×10^{-6} . Thus, the untyped SNPs from the low density
333 chip were less reliable than the newly discovered SNPs.

334 Genotyping Error Rate

335 Among 30 individuals (resulting in 65 pairwise comparisons, due to some birds being typed
336 >2 times) that were repeat genotyped on the SNP chip, there was a per SNP genotyping error
337 rate of 0.004. If comparisons were restricted to the 56 comparisons where both samples had
338 call rates >0.98, the error rate was 0.002, indicating that individuals with lower call rates
339 tended to be more error prone. The discordance in SNP calls between the chip and the
340 resequenced data was ~0.01, although this was apparently mostly driven by errors in the
341 sequencing data, because the degree of discordance is negatively correlated with the depth
342 of the genome coverage, which varies between 4.5x and 13.8x (see Fig. S3).

343

344 Resequencing data predict SNP chip allele frequencies

345 The minor allele frequencies (MAFs) of each SNP estimated from the 30 resequenced birds
 346 were compared to the MAFs estimated from the 996 birds genotyped in the Wytham Woods
 347 population. Notably, there was a very strong positive relationship between the minor allele
 348 frequencies in the two datasets (Fig. S4A; HD Chip MAF = $0.016 + 0.918 \cdot \text{ReSeq MAF}$, $F_{1,480756}$
 349 = $1.65 \cdot 10^6$, $r^2 = 0.77$, $P < 2.2 \cdot 10^{-16}$). Thus, the MAFs estimated from the resequencing data
 350 from 30 birds sampled across Europe are a reliable predictor of the MAFs obtained by typing
 351 a much larger sample from a single population on the HD chip. Similar analyses using
 352 genotyped birds from two randomly selected mainland European populations showed the
 353 same pattern (Fig. S4B, S4C); Montpellier, HD Chip MAF = $0.023 + 0.867 \cdot \text{ReSeq MAF}$, $F_{1,480756}$
 354 = $8.16 \cdot 10^5$, $r^2 = 0.63$, $P < 2.2 \cdot 10^{-16}$, 50 individuals; Gotland, HD Chip MAF = $0.022 +$
 355 $0.874 \cdot \text{ReSeq MAF}$, $F_{1,480756} = 8.69 \cdot 10^5$, $r^2 = 0.64$, $P < 2.2 \cdot 10^{-16}$, 47 individuals. The relationship
 356 was stronger for the Wytham Woods birds than the two other populations, but this is largely
 357 because the HD chip MAFs were estimated from more birds in the Wytham Woods dataset,
 358 and are therefore presumably estimated more accurately. A similar analysis conducted on 50
 359 randomly chosen birds from Wytham Woods produced a relationship that was only slightly
 360 stronger than that seen in the Montpellier and Gotland populations (Fig. S4D; HD Chip MAF =
 361 $0.023 + 0.879 \cdot \text{ReSeq MAF}$, $F_{1,480756} = 9.76 \cdot 10^5$, $r^2 = 0.67$, $P < 2.2 \cdot 10^{-16}$). Thus, the strong
 362 relationship between SNP chip MAF and resequencing is not simply an artefact of 10 of the
 363 30 resequenced birds being from Wytham Woods. The mean minor allele frequencies were
 364 very similar in the three populations (Wytham 0.280, Montpellier 0.273, Gotland 0.274).

365 CNV analysis

366 A total of 41 526 putative CNVs (34,947 with PennCNV confidence scores >5) were
 367 discovered in 996 birds from Wytham Woods. The great majority (37 419 or 90.1%) of CNVs
 368 were single copy duplications. Birds had a mean (SD) of 41.9 (160.9) CNVs each, spanning a
 369 mean (SD) distance of 3.19 (16.22) Mbp. However, there was a strong positive relationship
 370 between the amount of CNV in a bird's genome and the Axiom Analysis Suite parameter
 371 cluster_distance_SD (Figure 2). Cluster_distance_SD is a per-sample measure, defined as the
 372 standard deviation of the distance to the cluster centre, estimated from all of the individual's
 373 called genotypes. Samples with high values of cluster_distance_SD are typically indicative of
 374 individuals whose genotypes are difficult to call, perhaps because the sample was of low
 375 quality or quantity. Restricting the analysis to those individuals with cluster_distance_SD

376 <0.65 (n = 701), resulted in far fewer CNVs. In total there were 8139 CNVs observed, of which
377 1523 (18.7%) were a deletion of two copies (i.e. the segment was missing from both
378 chromosomes), 1 424 (17.5%) were single copy deletions, 5,176 were single copy
379 duplications (63.6%) and 16 (0.2%) were double copy duplications. The retained birds had a
380 mean (SD) of 11.6 (6.8) CNVs spanning a mean (SD) total distance of 0.34 (0.40) Mbp. The
381 distributions of the number and total distance spanned of CNVs in the full dataset were far
382 more skewed (Fig. 3A, 3B) than in the restricted dataset (Fig. 3C, 3D). The skewedness of the
383 number and total distance of CNVs in the full dataset was 10.98 and 10.93 respectively, while
384 equivalent values in the restricted dataset were 2.71 and 5.14. Plink estimated there were
385 1397 distinct non-overlapping CNVs, of which 1204 were at a frequency < 0.01. However, a
386 small number of CNVs were at a frequency approaching 0.15. For an example of a large CNV
387 identified in multiple individuals see Fig. S5.

388 PennCNV analysis of two replicates of the reference genome bird, Abel, revealed there were
389 fewer CNVs than in the Wytham Woods population. For one replicate, the
390 cluster_distance_SD score was sufficiently low (0.59) to retain the sample in the filtered
391 dataset. No CNVs were detected, which is perhaps not surprising as CNV regions may not
392 have been possible to assemble when the genome was being assembled. The other replicate
393 had a cluster_distance_SD score of 0.69, and contained a total of six possible CNVs (although
394 four of them had confidence scores <5), with a total length of 104 Kbp. Some, perhaps all, of
395 these CNVs are likely to be false positives, but even with their inclusion, the reference bird
396 contains less CNV regions than the mean of the Wytham Woods dataset (mean summed
397 CNVs = 3.19 Mbp in the unfiltered dataset, 0.34 Mbp in the filtered dataset).

398

399 An analysis of nine father-mother-offspring trios from Wytham Woods (Table S4) identified
400 103 possible CNVs, of which 98 showed Mendelian inheritance, suggesting they were likely to
401 be correct calls. 71 CNVs involved insertions, 27 involved deletions and 5 had both insertions
402 and deletions segregating at the same location. The ratio of insertions: deletions is similar to
403 that described in the analyses of all Wytham Woods samples.

404

405 Genetic architecture of Exploration Behaviour

406 The GWAS of EB did not identify any SNPs that were significant at the genome-wide level
407 (Figure 4A). The QQ plots indicated that the distribution of p values was very close to that
408 expected under the null distribution if none of the SNPs explain variation in EB (Figure 4B),
409 and lambda was estimated as 1.018 (SE 1.7×10^{-5}). Thus the effects of population genetic
410 structure seem to be adequately accounted for. However, one SNP approached genome-
411 wide significance ($P = 0.136$; Table S5), and is worthy of mention. SNP AX-100303447 at
412 49.67Mbp on Chromosome 3 is located approximately 3.5 Kbp downstream of interleukin 22
413 receptor subunit alpha 2 *IL22RA2* (Figure 4C). This gene is notable for being implicated in the
414 regulation of alcohol drinking in alcohol-preferring laboratory rats; experimental interference
415 of *IL22RA2* expression results in reduced alcohol intake (Franklin *et al.* 2015). There is no
416 evidence that the *DRD4* gene explains variation in exploration behaviour in the Wytham
417 Woods population (Figure 4D).

418 The BayesR analysis of EB was consistent with a highly polygenic genetic architecture. The
419 heritability estimate was modest and had a very large 95% credible interval (Table 3),
420 although it was very similar to previous estimates from pedigree-based quantitative genetic
421 analyses. It was estimated that a large number of SNPs contributed to trait variation, and that
422 much of the additive genetic variance (V_A) was caused by SNPs in the smaller effect size
423 distributions (Table 3).

424

425 Discussion

426 In this study, we generated a high density SNP chip and showed that the majority of target
427 SNPs could be genotyped reliably and accurately and across multiple great tit populations. A
428 total of approximately 900 million SNP genotypes were generated with considerably less than
429 1% typing error. Similar chips are routinely used in studies of humans (Frazer *et al.* 2007;
430 Simonson *et al.* 2010), model organisms (Yang *et al.* 2009), companion animals (Hayward *et al.*
431 *al.* 2016) and agriculturally important species (Rincon *et al.* 2011; Winfield *et al.* 2016), but
432 their application in wild vertebrate populations remains rare – although there are some
433 examples using 40-50K SNP chips, e.g. in Soay sheep (Johnston *et al.* 2013), collared
434 flycatchers (Kawakami *et al.* 2014; Silva *et al.* 2017) and house sparrows (Silva *et al.* 2017).
435 We found the cost of genotyping to be relatively low (approximately £0.0003 per SNP
436 genotype per individual).

437 Several lessons were learned that may be useful to researchers considering designing their
438 own HD chips. First, we attempted to type some samples that were of marginal quality
439 relative to the manufacturer's recommendations. Although many of them were successfully
440 typed, the pass rate was lower than the remaining samples. Second, our chip included some
441 SNPs that had already been successfully typed on a smaller 10K Illumina SNP chip. These SNPs
442 did perform better than those which were unproven prior to the HD chip manufacture. Thus,
443 we recommend using SNPs that have been previously validated, even if prior testing was
444 performed on an alternative platform. Third, in addition to sequencing 10 birds from Wytham
445 Woods, we sequenced 20 birds from multiple other populations during the SNP discovery
446 and there is little evidence that the chip is biased towards SNPs that are more polymorphic in
447 the Wytham Woods population. If the discovery had relied on sequencing a single population
448 it is likely that there would have been a greater ascertainment bias towards SNPs that have
449 high minor allele frequencies in that population. Perhaps, most importantly, our relatively
450 high success rate (~82% of attempted assays were converted to QC-passed, polymorphic
451 SNPs) is at least partially attributable to performing SNP calling with different datasets and
452 different callers and then using consensus SNPs for the chip design.

453 There was a strong positive correlation between the SNP MAFs predicted from the 30
454 resequenced birds during the discovery phase, and the chip MAFs estimated from almost

455 1000 genotyped birds from the Wytham Woods population. During SNP discovery there will
456 be a tendency to assign higher confidence scores to SNPs with higher MAFs, because the rare
457 allele will be identified in multiple individuals. Thus, the site frequency spectrum of the SNP
458 chip cannot be expected to be representative of the whole genome, but for many
459 applications, a chip with relatively high MAFs can be beneficial. This is most obviously the
460 case in GWAS or linkage mapping studies where the power to detect linkage is partially a
461 function of MAF. The chip has already been used to detect regions of the genome
462 responsible for adaptive evolution of bill length in European great tits (Bosse *et al.* 2017).

463

464 We used the chip to examine the genetic architecture of Exploration Behaviour, a widely-
465 studied behavioural trait in great tits (Fidler *et al.* 2007; Korsten *et al.* 2010; Mueller *et al.*
466 2013) and other bird species (Edwards *et al.* 2015). No SNP reached genome-wide
467 significance, although this is perhaps unsurprising given that the sample size was fairly
468 modest (~400) and the trait was shown to have a reasonably low heritability in this dataset.
469 These findings are similar to a previous study using the lower density 10K chip, where
470 heritability of EB was also modest ($h^2 = 0.26$, $SE = 0.08$) and no SNPs were significant at the
471 genome-wide level in a GWAS (Santure *et al.* 2015).

472 Previous candidate gene studies of personality traits in great tits and other birds have
473 focused mainly on dopamine receptor D4 (*DRD4*), and there is convincing evidence that it
474 explains a small but significant amount of variation in great tit EB in a population in the
475 Netherlands (Fidler *et al.* 2007). With this in mind, *DRD4* was chosen as a candidate gene
476 during the SNP construction and the region was over-represented on the chip. However,
477 there was very little evidence that *DRD4* explained significant variation in the Wytham Woods
478 population. This is consistent with earlier studies (Korsten *et al.* 2010; Mueller *et al.* 2013)
479 that failed to find an association in Wytham Woods and elsewhere. It is probably prudent to
480 be cautious about most associations between *DRD4* and exploration behaviour in bird
481 species, unless genome-wide data are available. This is because single locus studies are
482 unable to reveal the extent to which test statistic inflation due to population structure or
483 covariance between environmental and additive genetic variance is driving false positive
484 results; see for example Knowler *et al.* (1988), discussed in Lynch & Walsh (1998). Of course,

485 this potential form of bias applies to any candidate gene study that lacks comparable data
486 from numerous non-candidate genomic regions.

487

488 High density SNP chips have been used to identify structural or copy-number variation (CNVs)
489 in other organisms (Wang *et al.* 2013; Wu *et al.* 2015; Zhang *et al.* 2014). We used the
490 PennCNV software to identify putative CNVs in the great tit genome. CNVs tended to be at
491 low frequency, which made validation hard because relatively few cases of each putative CNV
492 are present. Furthermore, it was clear that lower quality samples were prone to false positive
493 CNV calls. An additional complexity is that identifying the exact start and end points of each
494 CNV is non-trivial, so when CNVs in different birds partially overlap, it is not straightforward
495 to determine whether they are the same CNV or not. That CNVs have lower minor allele
496 frequencies than SNPs is not surprising because (i) they may be under stronger purifying
497 selection if they have bigger phenotypic effects and (ii) the chip was biased in favour of the
498 inclusion of SNPs with moderately high minor allele frequencies and designed completely
499 blind to the existence of CNVs. While CNVs are not a main focus of this study, it is clear that
500 some CNV calls were repeatable across different birds, and that the extent and effects on
501 phenotypic variation of CNVs are legitimate follow-up questions. Future CNV analyses should
502 ideally include replication from different methodologies (e.g. qPCR or sequencing-based
503 methods).

504

505 High density chips provide a straightforward method for typing several hundred thousand
506 SNPs. It is also the case that HD chips are relatively robust to low yield or highly degraded
507 DNA, whereas the DNA requirements for sequencing, especially long-read sequencing
508 technologies, tend to be more demanding. Whole genome sequencing remains more
509 expensive than SNP typing on a per individual basis, but that will not be the case for much
510 longer. Indeed, the HD chip era may be relatively short. Sequencing strategies that involve
511 sequencing a few individuals' genomes at high coverage, which are then used to impute the
512 genomes of many more individuals sequenced at ~1x coverage or lower, may already be as
513 cheap an alternative, and will yield more data (Gorjanc *et al.* 2015; Li *et al.* 2011; Pasaniuc *et al.*
514 *et al.* 2012). At present low coverage whole genome sequencing results in data that are harder

515 to process, although the challenges of low coverage assembly, SNP calling and imputation are
516 becoming more straightforward. Ecological genomics studies that use low-coverage
517 sequencing of many individuals are not yet common, but there are a few notable examples
518 e.g. a population genomic analysis of walking-stick insects *Timema* genomes (Soria-Carrasco
519 *et al.* 2014) and a phylogeography study of *Menidia menidia*, the Atlantic silverside fish,
520 (Therkildsen & Palumbi 2017)

521 In summary, high density SNP chips are a relatively straightforward approach for investigating
522 a diverse range of evolutionary genomics topics such as genetic architecture, adaptive
523 evolution, phylogeography, and inbreeding depression. Ultimately HD chips will be replaced
524 by whole genome sequencing, but they are likely to be used for a few more years, especially
525 in population genetic studies of organisms with very large genome sizes such as pines (Neale
526 *et al.* 2014; Nystedt *et al.* 2013) or salamanders (Nowoshilow *et al.* 2018), where sequencing
527 remains a relatively expensive option. We hope that the methodologies, lessons learned and
528 downstream applications described in this paper will be useful to other researchers
529 considering developing a similar chip to address evolutionary or ecological questions in their
530 favourite study organism. The chip described in this paper is available to other users from
531 Thermo Fisher Scientific (the company that acquired Affymetrix in 2016). In great tits, the
532 chip has already been used to detect signatures of selection (Bosse *et al.* 2017), to perform
533 genomewide association studies on morphological (Bosse *et al.* 2017) and phenological
534 (Gienapp *et al.* 2017b) traits, and to carry out detailed analysis of the role of CNVs on
535 genomic architecture (da Silva *et al.* In Press).

536 **References**

- 537 Albrechtsen A, Nielsen FC, Nielsen R (2010) Ascertainment Biases in SNP Chips Affect Measures of
538 Population Divergence. *Molecular Biology and Evolution* **27**, 2534-2547.
- 539 Aplin LM, Farine DR, Morand-Ferron J, *et al.* (2015) Experimentally induced innovations lead to
540 persistent culture via conformity in wild birds. *Nature* **518**, 538-541.
- 541 Aplin LM, Farine DR, Morand-Ferron J, Sheldon BC (2012) Social networks predict patch discovery in a
542 wild population of songbirds. *Proceedings Of The Royal Society B-Biological Sciences* **279**,
543 4199-4205.
- 544 Aulchenko YS, de Koning D-J, Haley C (2007a) Genomewide rapid association using mixed model and
545 regression: A fast and simple method for genomewide pedigree-based quantitative trait loci
546 association analysis. *Genetics* **177**, 577-585.
- 547 Aulchenko YS, Ripke S, Isaacs A, Van Duijn CM (2007b) GenABEL: an R library for genome-wide
548 association analysis. *Bioinformatics* **23**, 1294-1296.
- 549 Bajgain P, Rouse MN, Anderson JA (2016) Comparing Genotyping-by-Sequencing and Single
550 Nucleotide Polymorphism Chip Genotyping for Quantitative Trait Loci Mapping in Wheat.
551 *Crop Science* **56**, 232-248.
- 552 Berenos C, Ellis PA, Pilkington JG, Pemberton JM (2014) Estimating quantitative genetic parameters in
553 wild populations: a comparison of pedigree and genomic approaches. *Molecular Ecology* **23**,
554 3434-3451.
- 555 Bosse M, Spurgin LG, Laine VN, *et al.* (2017) Recent natural selection causes adaptive evolution of an
556 avian polygenic trait. *Science* **358**, 365-368.
- 557 Brawand D, Wagner CE, Li YI, *et al.* (2014) The genomic substrate for adaptive radiation in African
558 cichlid fish. *Nature* **513**, 375-381.
- 559 Bruford MW, Hanotte O, Brookfield JFY, Burke T (1998) Multilocus and single-locus DNA
560 fingerprinting. In: *Molecular genetic analysis of populations: A practical approach* (ed. Hoelzel
561 AR), pp. 225-269. IRL Press, Oxford, U.K.
- 562 Chang C, Chow C, Tellier L, *et al.* (2015) Second-generation PLINK: rising to the challenge of larger and
563 richer datasets. *GigaScience* **4**, 7.
- 564 Charmantier A, McCleery RH, Cole LR, *et al.* (2008) Adaptive Phenotypic Plasticity in Response to
565 Climate Change in a Wild Bird Population. *Science* **320**, 800-803.
- 566 Colbourne JK, Pfrender ME, Gilbert D, *et al.* (2011) The ecoresponsive genome of *Daphnia pulex*.
567 *Science* **331**, 555-561.
- 568 Cole EF, Morand-Ferron J, Hinks AE, Quinn JL (2012) Cognitive Ability Influences Reproductive Life
569 History Variation in the Wild. *Current Biology* **22**, 1808-1812.
- 570 Cole EF, Quinn JL (2014) Shy birds play it safe: personality in captivity predicts risk responsiveness
571 during reproduction in the wild. *Biology Letters* **10**, 20140178.
- 572 Corcoran P, Gossmann TI, Barton HJ, Slate J, Zeng K (2017) Determinants of the efficacy of natural
573 selection on coding and noncoding variability in two passerine species. *Genome Biology and
574 Evolution*, 2987-3007.
- 575 Danecek P, Auton A, Abecasis G, *et al.* (2011) The variant call format and vcf tools. *Bioinformatics*. **27**,
576 2156-2158.
- 577 Davey JW, Hohenlohe PA, Etter PD, *et al.* (2011) Genome-wide genetic marker discovery and
578 genotyping using next-generation sequencing. *Nature Reviews Genetics* **12**, 499-510.
- 579 DePristo MA, Banks E, Poplin R, *et al.* (2011) A framework for variation discovery and genotyping
580 using next-generation DNA sequencing data. *Nature Genetics* **43**, 491-498.
- 581 Dingemanse NJ, Both C, Drent PJ, Tinbergen JM (2004) Fitness consequences of avian personalities in
582 a fluctuating environment. *Proceedings Of The Royal Society B-Biological Sciences* **271**, 847-
583 852.
- 584 Drent PJ, van Oers K, van Noordwijk AJ (2003) Realized heritability of personalities in the great tit
585 (*Parus major*). *Proceedings Of The Royal Society B-Biological Sciences* **270**, 45-51.

- 586 Ebstein RP, Novick O, Umansky R, *et al.* (1996) Dopamine D4 receptor (D4DR) exon III polymorphism
587 associated with the human personality trait of novelty seeking. *Nature Genetics* **12**, 78-80.
- 588 Edwards HA, Hajduk GK, Durieux G, Burke T, Dugdale HL (2015) No Association between Personality
589 and Candidate Gene Polymorphisms in a Wild Bird Population. *Plos One* **10**, e0138439.
- 590 Ellegren H, Smeds L, Burri R, *et al.* (2012) The genomic landscape of species divergence in *Ficedula*
591 flycatchers. *Nature* **491**, 756-760.
- 592 Erbe M, Hayes BJ, Matukumalli LK, *et al.* (2012) Improving accuracy of genomic predictions within and
593 between dairy cattle breeds with imputed high-density single nucleotide polymorphism
594 panels. *Journal Of Dairy Science* **95**, 4114-4129.
- 595 Fidler AE, van Oers K, Drent PJ, *et al.* (2007) *Drd4* gene polymorphisms are associated with personality
596 variation in a passerine bird. *Proceedings of the Royal Society B: Biological Sciences* **274**, 1685-
597 1691.
- 598 Franklin KM, Hauser SR, Lasek AW, *et al.* (2015) Reduction of alcohol drinking of alcohol-preferring (P)
599 and high-alcohol drinking (HAD1) rats by targeting phosphodiesterase-4 (PDE4).
600 *Psychopharmacology* **232**, 2251-2262.
- 601 Frazer KA, Ballinger DG, Cox DR, *et al.* (2007) A second generation human haplotype map of over 3.1
602 million SNPs. *Nature* **449**, 851-853.
- 603 Garant D, Kruuk LEB, Wilkin TA, McCleery RH, Sheldon BC (2005) Evolution driven by differential
604 dispersal within a wild bird population. *Nature* **433**, 60-65.
- 605 Garroway CJ, Radersma R, Sepil I, *et al.* (2013) Fine-scale genetic structure in a wild bird population:
606 the role of limited dispersal and environmentally based selection as causal factors. *Evolution*
607 **67**, 3488-3500.
- 608 Gienapp P, Fior S, Guillaume F, *et al.* (2017a) Genomic Quantitative Genetics to Study Evolution in the
609 Wild. *Trends in Ecology & Evolution* **32**, 897-908.
- 610 Gienapp P, Laine VN, Mateman AC, van Oers K, Visser ME (2017b) Environment-Dependent
611 Genotype-Phenotype Associations in Avian Breeding Time. *Front Genet* **8**, 102.
- 612 Gorjanc G, Cleveland MA, Houston RD, Hickey JM (2015) Potential of genotyping-by-sequencing for
613 genomic selection in livestock populations. *Genetics Selection Evolution* **47**, 12.
- 614 Gossman TI, Santure AW, Sheldon BC, Slate J, Zeng K (2014) Highly Variable Recombinational
615 Landscape Modulates Efficacy of Natural Selection in Birds. *Genome Biology and Evolution* **6**,
616 2061-2075.
- 617 Groothuis TGG, Carere C (2005) Avian personalities: characterization and epigenesis. *Neuroscience*
618 *and Biobehavioral Reviews* **29**, 137-150.
- 619 Harvey PH, Greenwood PJ, Perrins CM (1979) Breeding Area Fidelity of Great Tits (*Parus-Major*).
620 *Journal Of Animal Ecology* **48**, 305-313.
- 621 Hayward JJ, Castelhana MG, Oliveira KC, *et al.* (2016) Complex disease and phenotype mapping in the
622 domestic dog. *Nature Communications* **7**, 10460.
- 623 Hohenlohe PA, Bassham S, Etter PD, *et al.* (2010) Population Genomics of Parallel Adaptation in
624 Threespine Stickleback using Sequenced RAD Tags. *PLoS Genet* **6**, e1000862.
- 625 Hu TT, Pattyn P, Bakker EG, *et al.* (2011) The *Arabidopsis lyrata* genome sequence and the basis of
626 rapid genome size change. *Nature Genetics* **43**, 476-481.
- 627 Iamartino D, Nicolazzi EL, Van Tassell CP, *et al.* (2017) Design and validation of a 90K SNP genotyping
628 assay for the water buffalo (*Bubalus bubalis*). *Plos One* **12**.
- 629 Jensen H, Szulkin M, Slate J (2014) Molecular Quantitative Genetics. In: *Quantitative Genetics in the*
630 *Wild* (eds. Charmantier A, Garant D, Kruuk LEB), pp. 209-227. Oxford University Press, Oxford.
- 631 Johnston SE, Gratten J, Berenos C, *et al.* (2013) Life history trade-offs at a single locus maintain
632 sexually selected genetic variation. *Nature* **502**, 93-95.
- 633 Jones FC, Grabherr MG, Chan YF, *et al.* (2012) The genomic basis of adaptive evolution in threespine
634 sticklebacks. *Nature* **484**, 55-61.

- 635 Kawakami T, Backstrom N, Burri R, *et al.* (2014) Estimation of linkage disequilibrium and interspecific
636 gene flow in Ficedula flycatchers by a newly developed 50k single-nucleotide polymorphism
637 array. *Molecular Ecology Resources* **14**, 1248-1260.
- 638 Kluijver HN (1951) *The Population Ecology of the Great Tit, Parus M. Major* L Brill.
- 639 Knowler WC, Williams RC, Pettitt DJ, Steinberg AG (1988) Gm3-5,13,14 and Type-2 Diabetes-Mellitus -
640 an Association in American-Indians with Genetic Admixture. *American Journal of Human*
641 *Genetics* **43**, 520-526.
- 642 Korneliussen TS, Albrechtsen A, Nielsen R (2014) ANGSD: Analysis of Next Generation Sequencing
643 Data. *BMC Bioinformatics* **15**, 356.
- 644 Korsten P, Mueller JC, Hermannstadter C, *et al.* (2010) Association between DRD4 gene polymorphism
645 and personality variation in great tits: a test across four wild populations. *Molecular Ecology*
646 **19**, 832-843.
- 647 Kranis A, Gheyas AA, Boschiero C, *et al.* (2013) Development of a high density 600K SNP genotyping
648 array for chicken. *Bmc Genomics* **14**, 59.
- 649 Kvist L, Martens J, Higuchi H, *et al.* (2003) Evolution and genetic structure of the great tit (Parus
650 major) complex. *Proceedings Of The Royal Society B-Biological Sciences* **270**, 1447-1454.
- 651 Lack D (1964) A Long-Term Study of the Great Tit (Parus major). *Journal Of Animal Ecology* **33**, 159-
652 173.
- 653 Lack D (1968) *Ecological adaptations for breeding in birds* Methuen, London.
- 654 Laine VN, Gossman TI, Schachtschneider KM, *et al.* (2016) Evolutionary signals of selection on
655 cognition from the great tit genome and methylome. *Nature Communications* **7**, 10474.
- 656 Lamichhaney S, Berglund J, Almen MS, *et al.* (2015) Evolution of Darwin's finches and their beaks
657 revealed by genome sequencing. *Nature* **518**, 371-375.
- 658 Lemoine M, Lucek K, Perrier C, *et al.* (2016) Low but contrasting neutral genetic differentiation shaped
659 by winter temperature in European great tits. *Biological Journal of the Linnean Society* **118**,
660 668-685.
- 661 Li H (2011) A statistical framework for SNP calling, mutation discovery, association mapping and
662 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987-
663 2993.
- 664 Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform.
665 *Bioinformatics* **25**, 1754-1760.
- 666 Li H, Handsaker B, Wysoker A, *et al.* (2009) The Sequence Alignment/Map format and SAMtools.
667 *Bioinformatics* **25**, 2078-2079.
- 668 Li Y, Sidore C, Kang HM, Boehnke M, Abecasis GR (2011) Low-coverage sequencing: Implications for
669 design of complex trait association studies. *Genome Research* **21**, 940-951.
- 670 Liu S, Sun L, Li Y, *et al.* (2014) Development of the catfish 250K SNP array for genome-wide association
671 studies. *BMC Res Notes* **7**, 135.
- 672 Lynch M, Walsh B (1998) *Genetics and Analysis of Quantitative Traits* Sinauer, Sunderland,
673 Massachusetts.
- 674 Matsuzaki H, Dong SL, Loi H, *et al.* (2004) Genotyping over 100,000 SNPs on a pair of oligonucleotide
675 arrays. *Nature Methods* **1**, 109-111.
- 676 McCarroll SA, Altshuler DM (2007) Copy-number variation and association studies of human disease.
677 *Nature Genetics* **39**, S37-S42.
- 678 McCleery RH, Pettifor RA, Armbruster P, *et al.* (2004) Components of variance underlying fitness in a
679 natural population of the great tit Parus major. *American Naturalist* **164**, E62-E72.
- 680 McKenna A, Hanna M, Banks E, *et al.* (2010) The Genome Analysis Toolkit: A MapReduce framework
681 for analyzing next-generation DNA sequencing data. *Genome Research* **20**, 1297-1303.
- 682 Miller JM, Kijas JW, Heaton MP, McEwan JC, Coltman DW (2012) Consistent divergence times and
683 allele sharing measured from cross-species application of SNP chips developed for three
684 domestic species. *Molecular Ecology Resources* **12**, 1145-1150.

- 685 Moser G, Lee SH, Hayes BJ, *et al.* (2015) Simultaneous Discovery, Estimation and Prediction Analysis of
686 Complex Traits Using a Bayesian Mixture Model. *Plos Genetics* **11**, e1004969.
- 687 Mueller JC, Korsten P, Hermannstaedter C, *et al.* (2013) Haplotype structure, adaptive history and
688 associations with exploratory behaviour of the DRD4 gene region in four great tit (*Parus*
689 major) populations. *Molecular Ecology* **22**, 2797-2809.
- 690 Nagamine Y, Pong-Wong R, Navarro P, *et al.* (2012) Localising Loci underlying Complex Trait Variation
691 Using Regional Genomic Relationship Mapping. *Plos One* **7**, 12.
- 692 Neale DB, Wegrzyn JL, Stevens KA, *et al.* (2014) Decoding the massive genome of loblolly pine using
693 haploid DNA and novel assembly strategies. *Genome Biology* **15**, R59.
- 694 Nowoshilow S, Schloissnig S, Fei JF, *et al.* (2018) The axolotl genome and the evolution of key tissue
695 formation regulators. *Nature* **554**, 50-55.
- 696 Nussey DH, Postma E, Gienapp P, Visser ME (2005) Selection on heritable phenotypic plasticity in a
697 wild bird population. *Science* **310**, 304-306.
- 698 Nystedt B, Street NR, Wetterbom A, *et al.* (2013) The Norway spruce genome sequence and conifer
699 genome evolution. *Nature* **497**, 579-584.
- 700 Pasaniuc B, Rohland N, McLaren PJ, *et al.* (2012) Extremely low-coverage sequencing and imputation
701 increases power for genome-wide association studies. *Nature Genetics* **44**, 631-641.
- 702 Perrins CM (1979) *British Tits* Collins.
- 703 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE (2012) Double Digest RADseq: An Inexpensive
704 Method for De Novo SNP Discovery and Genotyping in Model and Non-Model Species. *Plos*
705 *One* **7**, e37135.
- 706 Postma E, van Noordwijk AJ (2005) Gene flow maintains a large genetic difference in clutch size at a
707 small spatial scale. *Nature* **433**, 65-68.
- 708 Powell JE, Visscher PM, Goddard ME (2010) Reconciling the analysis of IBD and IBS in complex trait
709 studies. *Nature Reviews Genetics* **11**, 800-805.
- 710 Quinn JL, Cole EF, Reed TE, Morand-Ferron J (2016) Environmental and genetic determinants of
711 innovativeness in a natural population of birds. *Philosophical Transactions Of The Royal*
712 *Society B-Biological Sciences* **371**, 20150184.
- 713 Quinn JL, Patrick SC, Bouwhuis S, Wilkin TA, Sheldon BC (2009) Heterogeneous selection on a
714 heritable temperament trait in a variable environment. *Journal Of Animal Ecology* **78**, 1203-
715 1215.
- 716 Radersma R, Garroway CJ, Santure AW, *et al.* (2017) Social and spatial effects on genetic variation
717 between foraging flocks in a wild bird population. *Molecular Ecology* **26**, 5807-5819.
- 718 Rincon G, Weber KL, Van Eenennaam AL, Golden BL, Medrano JF (2011) Hot topic: Performance of
719 bovine high-density genotyping platforms in Holsteins and Jerseys. *Journal Of Dairy Science*
720 **94**, 6116-6121.
- 721 Robinson MR, Santure AW, DeCauwer I, Sheldon BC, Slate J (2013) Partitioning of genetic variation
722 across the genome using multimarker methods in a wild bird population. *Molecular Ecology*
723 **22**, 3963-3980.
- 724 Robledo D, Palaikostas C, Bargelloni L, Martínez P, Houston R (2017) Applications of genotyping by
725 sequencing in aquaculture breeding and genetics. *Reviews in Aquaculture*,
726 doi:10.1111/raq.12193.
- 727 Santure AW, De Cauwer I, Robinson MR, *et al.* (2013) Genomic dissection of variation in clutch size
728 and egg mass in a wild great tit (*Parus major*) population. *Molecular Ecology* **22**, 3949-3962.
- 729 Santure AW, Gratten J, Mossman JA, Sheldon BC, Slate J (2011) Characterisation of the transcriptome
730 of a wild great tit *Parus major* population by next generation sequencing. *Bmc Genomics* **12**.
- 731 Santure AW, Poissant J, De Cauwer I, *et al.* (2015) Replicated analysis of the genetic architecture of
732 quantitative traits in two wild great tit populations. *Molecular Ecology* **24**, 6148-6162.
- 733 Sepil I, Lachish S, Hinks AE, Sheldon BC (2013) Mhc supertypes confer both qualitative and
734 quantitative resistance to avian malaria infections in a wild bird population. *Proceedings Of*
735 *The Royal Society B-Biological Sciences* **280**, 20130134.

- 736 Sepil I, Moghadam HK, Huchard E, Sheldon BC (2012) Characterization and 454 pyrosequencing of
737 Major Histocompatibility Complex class I genes in the great tit reveal complexity in a
738 passerine system. *Bmc Evolutionary Biology* **12**.
- 739 Shen R, Fan JB, Campbell D, *et al.* (2005) High-throughput SNP genotyping on universal bead arrays.
740 *Mutation Research-Fundamental and Molecular Mechanisms of Mutagenesis* **573**, 70-82.
- 741 Silva CNS, McFarlane SE, Hagen IJ, *et al.* (2017) Insights into the genetic architecture of morphological
742 traits in two passerine bird species. *Heredity* **119**, 197-205.
- 743 da Silva VH, Laine VN, Bosse M, van Oers K, Dibbits B, Visser ME, Crooijmans RPMA, Groenen MAM.
744 CNVs are associated with genomic architecture in a songbird. *BMC Genomics*, In Press.
- 745 Simonson TS, Yang YZ, Huff CD, *et al.* (2010) Genetic Evidence for High-Altitude Adaptation in Tibet.
746 *Science* **329**, 72-75.
- 747 Smith HG, Kallander H, Nilsson JA (1989) The Trade-Off between Offspring Number and Quality in the
748 Great Tit *Parus-Major*. *Journal Of Animal Ecology* **58**, 383-401.
- 749 Soria-Carrasco V, Gompert Z, Comeault AA, *et al.* (2014) Stick Insect Genomes Reveal Natural
750 Selection's Role in Parallel Speciation. *Science* **344**, 738-742.
- 751 Spencer CCA, Su Z, Donnelly P, Marchini J (2009) Designing Genome-Wide Association Studies: Sample
752 Size, Power, Imputation, and the Choice of Genotyping Chip. *Plos Genetics* **5**, e1000477.
- 753 Syvanen AC (2001) Accessing genetic variation: Genotyping single nucleotide polymorphisms. *Nature*
754 *Reviews Genetics* **2**, 930-942.
- 755 Therkildsen NO, Palumbi SR (2017) Practical low-coverage genomewide sequencing of hundreds of
756 individually barcoded samples for population and evolutionary genomics in nonmodel
757 species. *Mol Ecol Resour* **17**, 194-208.
- 758 Titulaer M, van Oers K, Naguib M (2012) Personality affects learning performance in difficult tasks in a
759 sex-dependent way. *Animal Behaviour* **83**, 723-730.
- 760 van Bers N, van Oers K, Kerstens H, *et al.* (2010) Genome-wide SNP detection in the great tit *Parus*
761 *major* using high throughput sequencing. *Molecular Ecology* **19**, 89-99.
- 762 Van Bers NEM, Santure AW, Van Oers K, *et al.* (2012) The design and cross-population application of a
763 genome-wide SNP chip for the great tit *Parus major*. *Molecular Ecology Resources* **12**, 753-
764 770.
- 765 van Oers K, de Jong G, Drent PJ, van Noordwijk AJ (2004) A genetic analysis of avian personality traits:
766 Correlated, response to artificial selection. *Behavior Genetics* **34**, 611-619.
- 767 Van Oers K, Naguib M (2013) Avian personality. In: *Animal Personalities: Behaviour, Physiology and*
768 *Evolution* (eds. Carere C, Maestriperi D), pp. 66-95. The University of Chicago Press, Chicago.
- 769 van Oers K, Santure AW, De Cauwer I, *et al.* (2014) Replicated high-density genetic maps of two great
770 tit populations reveal fine-scale genomic departures from sex-equal recombination rates.
771 *Heredity* **112**, 307-316.
- 772 Visscher PM, Wray NR, Zhang Q, *et al.* (2017) 10 Years of GWAS Discovery: Biology, Function, and
773 Translation. *American Journal of Human Genetics* **101**, 5-22.
- 774 Visser ME, Adriaensen F, van Balen JH, *et al.* (2003) Variable responses to large-scale climate change
775 in European *Parus* populations. *Proceedings of the Royal Society of London Series B-Biological*
776 *Sciences* **270**, 367-372.
- 777 Visser ME, van Noordwijk AJ, Tinbergen JM, Lessells CM (1998) Warmer springs lead to mistimed
778 reproduction in great tits (*Parus major*). *Proceedings of the Royal Society of London Series B-*
779 *Biological Sciences* **265**, 1867-1870.
- 780 Wang JY, Wang HF, Jiang JC, *et al.* (2013) Identification of Genome-Wide Copy Number Variations
781 among Diverse Pig Breeds Using SNP Genotyping Arrays. *Plos One* **8**, e68683.
- 782 Wang K, Li MY, Hadley D, *et al.* (2007) PennCNV: An integrated hidden Markov model designed for
783 high-resolution copy number variation detection in whole-genome SNP genotyping data.
784 *Genome Research* **17**, 1665-1674.
- 785 Winfield MO, Allen AM, Burrige AJ, *et al.* (2016) High-density SNP genotyping array for hexaploid
786 wheat and its secondary and tertiary gene pool. *Plant Biotechnology Journal* **14**, 1195-1206.

- 787 Wu Y, Fan HZ, Jing SY, *et al.* (2015) A genome-wide scan for copy number variations using high-density
788 single nucleotide polymorphism array in Simmental cattle. *Animal Genetics* **46**, 289-298.
789 Yang H, Ding YM, Hutchins LN, *et al.* (2009) A customized and versatile high-density genotyping array
790 for the mouse. *Nature Methods* **6**, 663-666.
791 Yang J, Manolio TA, Pasquale LR, *et al.* (2011) Genome partitioning of genetic variation for complex
792 traits using common SNPs. *Nature Genetics* **43**, 519-525.
793 Zhang X, Du RQ, Li SL, *et al.* (2014) Evaluation of copy number variation detection for a SNP array
794 platform. *BMC Bioinformatics* **15**, 50.

795

796 Author Contributions

797 Designed the chip (J-MK, AS, KvO, MAMG, JS), performed the CNV analysis (HB, JS), designed
798 and performed the exploration behaviour assays (JQ, EC), coordinated the long term data and
799 blood sample collection (MV, BS), performed the genetic architecture analyses (JS), collected
800 field and DNA sample data (Great Tit HapMap consortium), wrote the paper (JS, with
801 contributions from all authors), conceived the study (JS, BS, MAMG, KvO, MV).

802 Acknowledgements

803 This work was supported by grants from the European Research Council (grants 339092 to
804 M.E.V., 250164 to B.C.S., and 202487 to J.S.) and Natural Environment Research Council
805 (grant NE/J012599/1 to J.S.). We thank Claire Bloor, Geoff Scopes and Alessandro Davassi of
806 Affymetrix for their help during the chip design and genotyping calling processes. Richard
807 Talbot and Alison Downing of Edinburgh Genomics provided the genotyping service. Padraic
808 Corcoran, Kai Zeng, Veronika Laine, Lewis Spurgin and Mirte Bosse provided useful
809 discussions and or advice on SNP calling from Illumina reads or SNP chip quality control. We
810 thank Dany Garant and two anonymous reviewers for their insightful comments on, and
811 suggested improvements to, an earlier draft of the manuscript.

812

813 Data deposition

814 Genotype and phenotype data are deposited as Plink files on Dryad under the provisional
815 record doi:10.5061/dryad.7d467b6. All SNPs included on the chip are reported on the
816 European Variation Archive (<https://www.ebi.ac.uk/eva/>) under accession number
817 PRJEB24964. SNP discovery was performed on 30 resequenced birds, whose genomes are
818 reported on the NCBI sequence read archive under project ID SRP066678.

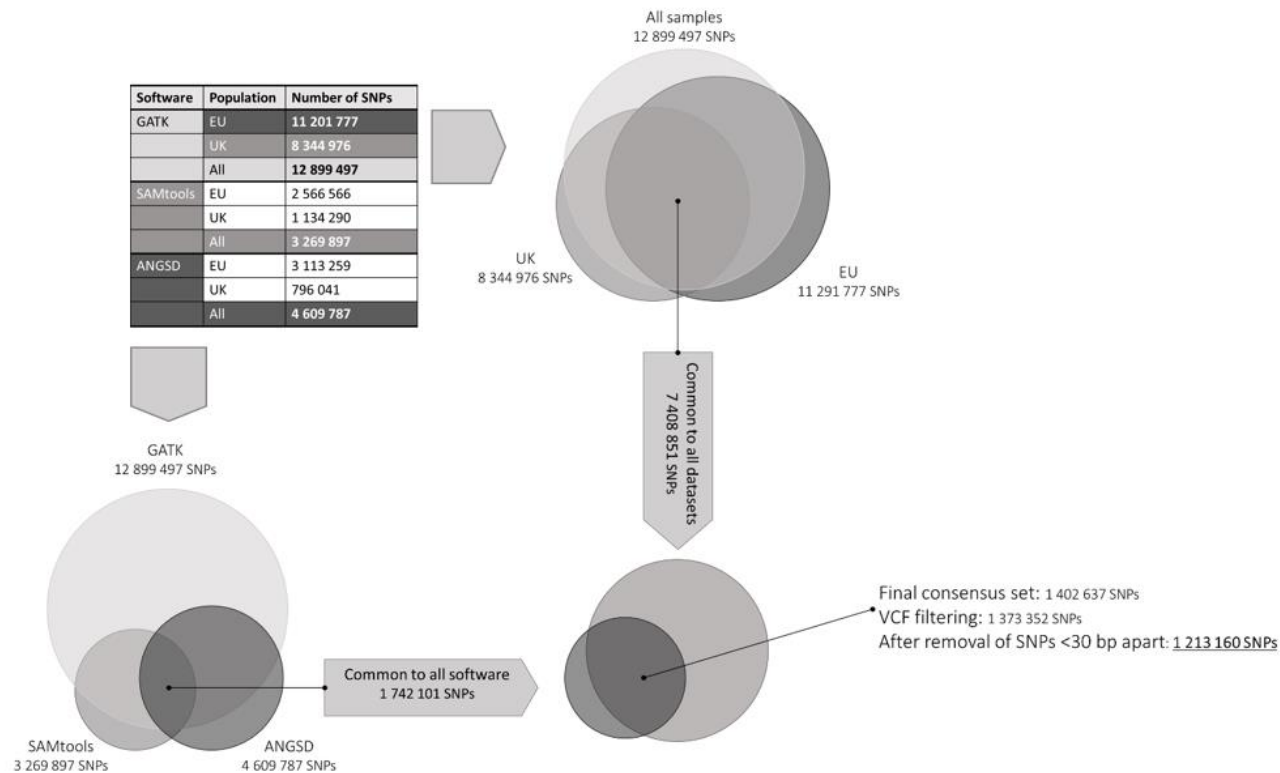
819

820 The Authors have no conflicts of interest.

821

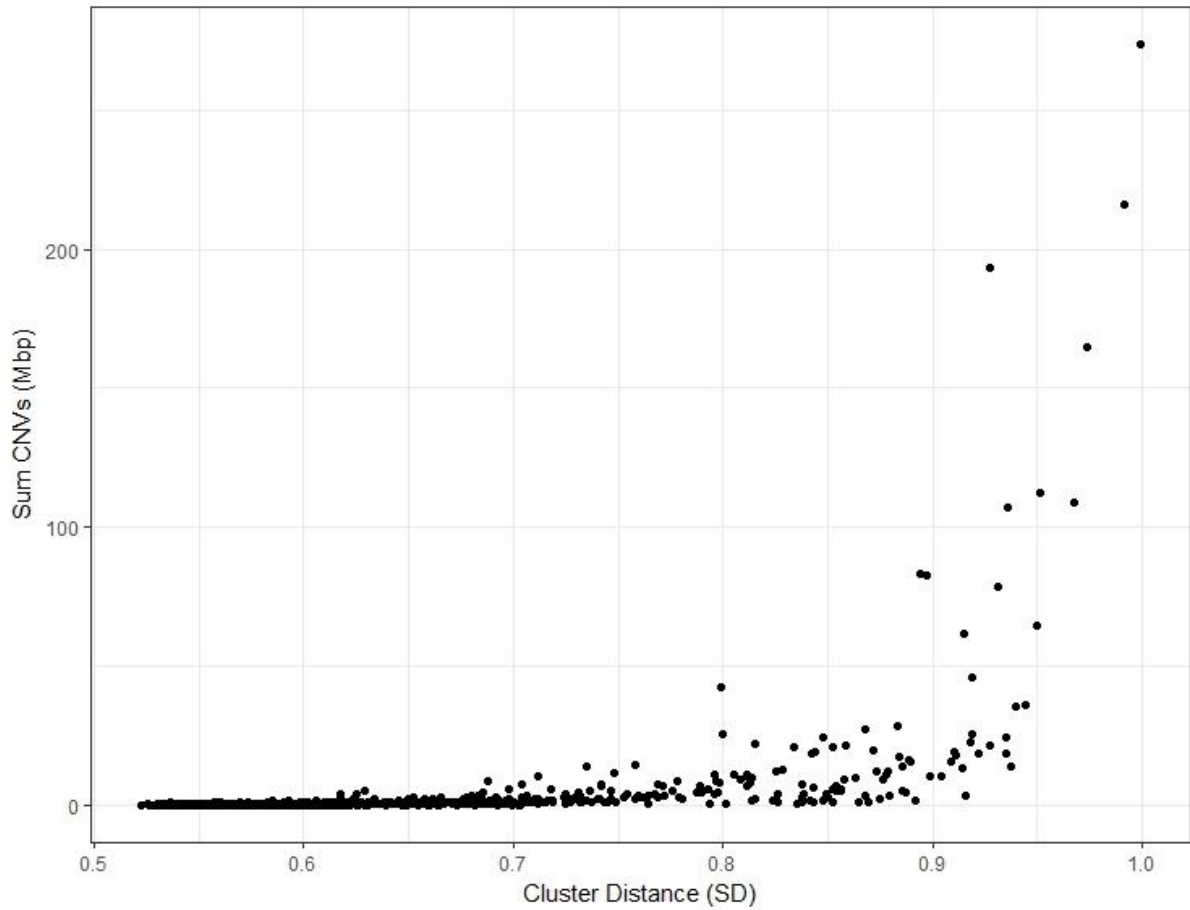
822 Figure 1: The pipeline for SNP discovery. The top right part of the figure identifies SNPs that were found in the UK birds (n= 10), the mainland Europe birds
 823 (n=20) and in all birds (n=30) with the software GATK. The bottom left part of the figure identified SNPs that were found when all 30 birds were analysed with
 824 three software packages - GATK, SAMtools and ANGSD. The intersection of these discovery pipelines, i.e. SNPs that were detected in all populations by all
 825 software packages, were considered for inclusion on the chip. After filtering for MAF > 0.05 and removal of SNPs located within 30bp of each other, a final list
 826 of 1 213 160 SNPs remained.

827



828 Figure 2: Individuals with higher standard deviation (SD) in their cluster distance, indicating samples
829 whose genotypes are difficult to call, tend to have a greater proportion of their genomes called as
830 CNVs. The assembled great tit genome is approximately 1020Mbp long.

831



832

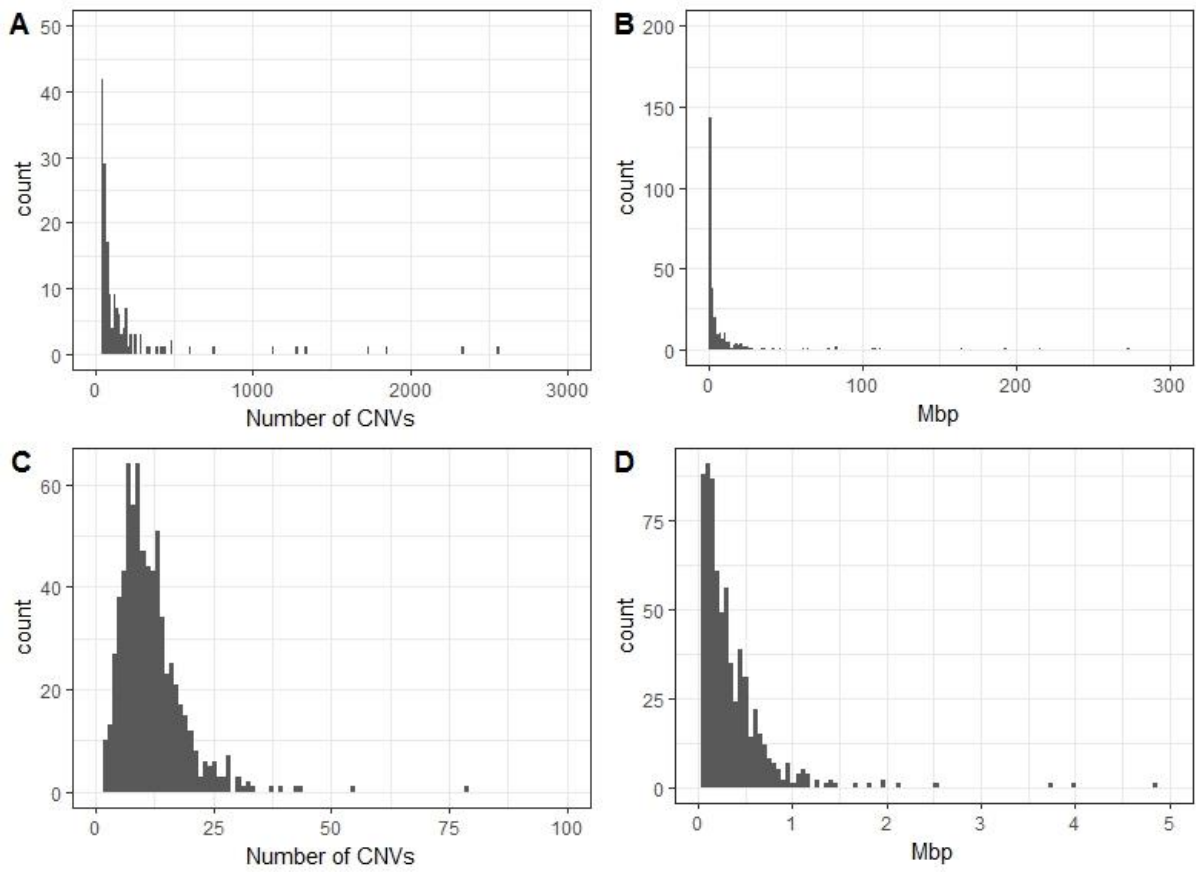
833

834

835 Figure 3: Distribution of the number and total distance spanned of CNVs in 996 Wytham Woods birds
836 (top panels) and the remaining 701 Wytham Woods birds after filtering on $\text{cluster_distance_SD} < 0.65$
837 (bottom panels); i.e. after removing samples whose genotypes are difficult to call.

838

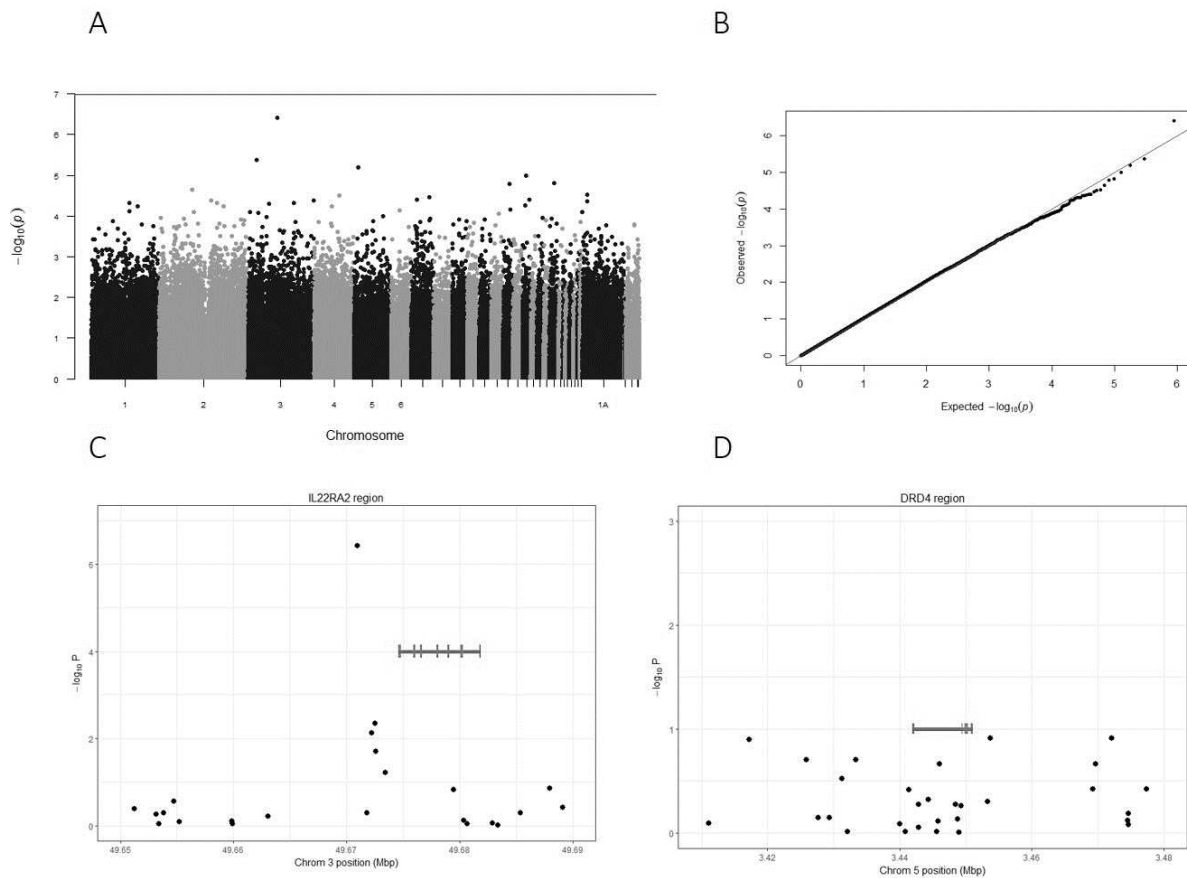
839



840

841 Figure 4: (A) Manhattan plot of a GWAS for Exploration Behaviour on 415 individuals for 459 502
 842 autosomal markers. Chromosomes are ordered numerically from 1-15, 17-24, 25LG1, 25LG2, 26-28,
 843 1A, 4A, LG22 and chromosome unknown. Horizontal line = genomewide significance. (B) QQ plot of
 844 observed versus expected $-\log_{10}$ transformed P values. Lambda = 1.018 (SE 1.7×10^{-7}). (C) and (D)
 845 Zoomed in plots of GWAS results close to the *IL22RA2* and *DRD4* genes. Horizontal lines represent
 846 location of genes. Note that the y-axis scale differs between plots.

847



848

849

850

851 Table 1: Great tit populations where genotyping was attempted (see also Fig. S1).

Population	Population Code	Coordinates (N, E) in decimal degrees	Birds typed	Birds passing QC
Amur, Russia ¹	1	50.62, 131.37	72	63
Antwerp, Belgium ³	2	51.13, 4.53	36	30
Cambridge, UK ³	3	52.40, -0.23	35	34
Font Roja, Spain ³	4	38.66, -0.54	30	29
Gotland, Sweden ³	5	57.14, 18.33	50	47
Groblas, Poland ³	6	52.28, 17.90	4	4
Harjavalta, Finland	7	61.33, 22.17	44	44
Hoge Veluwe, Netherlands	8	52.07, 5.84	38	36
Israel	9	32.62, 35.24	1	1
La Rouviere, France ³	10	43.66, 3.67	31	27
Loch Lomond, Scotland ³	11	56.13, -4.62	43	41
Mariola, Spain ³	12	38.73, -0.55	33	33
Montpellier, France ³	13	43.61, 3.87	50	50
Oulu, Finland ³	14	65.13, 25.88	50	45
Pilis Mountains, Hungary ³	15	47.72, 19.02	36	34
Pirio and Muro, Corsica ³	16	42.37, 8.75	30	27
Radolfzell, Germany	17	47.74, 8.98	30	27
Sakhalin Island, Russia ²	18	50.52, 143.11	13	13
Seewisen, Germany ³	19	47.97, 8.98	50	46
Tartu, Estonia ³	20	58.17, 25.08	43	42
Tomakomai, Japan ²	21	42.67, 141.60	10	9
Velky Kosir, Czech Republic ³	22	49.53, 17.07	36	33
Vienna, Austria ³	23	48.21, 16.26	38	31
Vlieland, Netherlands	24	53.28, 5.01	30	21
Westerheide, Netherlands	25	52.00, 5.83	39	35
Wytham Woods, UK ³	26	51.77, -1.33	1073	996
Zurich, Switzerland ³	27	47.39, 8.57	30	29
Zvenigorod, Russia	28	55.73, 36.85	20	19
Total			2007	1846

852 ¹ Sample contains 63 *Parus major* and 9 putative *P. major/P. minor* hybrids853 ² *Parus minor* populations854 ³ Population included in the 30 resequenced genomes dataset

855

856 Table 2: Summary of SNP genotype calling, by Affymetrix Axiom Analysis Suite category. The
 857 Conversion Type columns uses the Affymetrix terminology but can be summarised as follows:
 858 PolyHighResolution = SNP that is polymorphic and can be reliably scored due to the different
 859 genotypes forming resolvable, discrete clusters; NoMinorHom = similar to a PolyHighResolution, but
 860 where the minor allele homozygote is not observed, presumably due to a low genotype frequency;
 861 MonoHighResolution = a monomorphic SNP that can be reliably scored because it forms a single
 862 cluster; CallRateBelowThreshold = a SNP with the expected number of clusters (usually 3, one for
 863 each possible genotype), but where the proportion of samples scored at the SNP falls below a user-
 864 defined threshold. Here the threshold was 0.97; Off-target variant = SNPs, where additional (i.e. more
 865 than 3) clusters are observed, making genotype calling ambiguous; Other = all other unresolvable
 866 SNPs.

867

Conversion Type	Count	Percentage	Retained for analysis
PolyHighResolution	498 036	81.5	497 972
NoMinorHom	4048	0.7	4047
MonoHighResolution	666	0.1	666
CallRateBelowThreshold	40 499	6.6	0
Off Target Variant (OTV)	9545	1.6	0
Other	58 176	9.5	0
Sum	610 970		502 685

868

869

870

871 Table 3: Genetic architecture of exploration behaviour

Parameter	Estimate (95% credible interval)
Heritability	0.161 (<0.001-0.671)
Number of SNPs	3,253 (315-8,499)
PGE_0.0001	0.41 (0.01-0.89)
PGE_0.001	0.26 (<0.01-0.80)
PGE_0.01	0.33 (<0.01-0.90)

872 Heritability is the total heritability captured by the genotyped SNPs (often termed “SNP heritability” or
873 “chip heritability”). Number of SNPs is the number of SNPs inferred as explaining some (non-zero)
874 trait variation. PGE is the proportion of SNP heritability explained by SNPs in the 0.001, 0.001 and 0.01
875 effect size distributions.

876