

MOLECULAR ECOLOGY RESOURCES

The easy road to genome-wide medium density SNP screening in a non-model species: development and application of a 10K SNP-chip for the house sparrow (*Passer domesticus*).

Journal:	<i>Molecular Ecology Resources</i>
Manuscript ID:	MER-12-0405.R1
Manuscript Type:	Resource Article
Date Submitted by the Author:	11-Jan-2013
Complete List of Authors:	Hagen, Ingerid; Centre for Conservation Biology, Norwegian University of Science and Technology, Dept. of Biology Billing, Anna; Centre for Conservation Biology, Norwegian University of Science and Technology, Dept. of Biology Rønning, Bernt; Centre for Conservation Biology, Norwegian University of Science and Technology, Dept. of Biology Pedersen, Sindre; Norwegian University of Science and Technology, Dept. of Biology Pärn, Henrik; Centre for Conservation Biology, Norwegian University of Science and Technology, Dept. of Biology Slate, Jon; University of Sheffield, Dept. of Animal and Plant Sciences Jensen, Henrik; Centre for Conservation Biology, Norwegian University of Science and Technology, Dept. of Biology
Keywords:	next-generation sequencing, single nucleotide polymorphism, genetic population differentiation, house sparrow, <i>Passer domesticus</i>

1 **Title:**

2 **The easy road to genome-wide medium density SNP screening in a non-model species:**
3 **development and application of a 10K SNP-chip for the house sparrow (*Passer***
4 ***domesticus*).**

5

6 **Ingerid J. Hagen¹, Anna M. Billing¹, Bernt Rønning¹, Sindre A. Pedersen², Henrik**
7 **Pärn¹, Jon Slate³ and Henrik Jensen¹**

8

9 ¹ Centre for Conservation Biology, Department of Biology, Norwegian University of Science
10 and Technology, NO-7491 Trondheim, Norway.

11 ² Department of Biology, Norwegian University of Science and Technology, NO-7491
12 Trondheim, Norway.

13 ³ Department of Animal and Plant Sciences, University of Sheffield, Western Bank, Sheffield
14 S10 2TN, UK.

15

16 **Keywords:** next-generation sequencing, single nucleotide polymorphism, genetic population
17 differentiation, house sparrow, *Passer domesticus*

18

19 **Corresponding authors:** Ingerid J. Hagen, Fax: +47 73596100, E-mail:

20 ingerid.hagen@bio.ntnu.no, and Henrik Jensen, Fax: +47 73596100, E-mail:

21 henrik.jensen@bio.ntnu.no

22

23 **Running title:** Medium density SNP screening in sparrows

24

25 **Abstract**

26 With the advent of next generation sequencing, new avenues have opened to study genomics
27 in wild populations of non-model species. Here, we describe a successful approach to a
28 genome-wide medium density Single Nucleotide Polymorphism (SNP) panel in a non-model
29 species, the house sparrow (*Passer domesticus*), through the development of a 10K Illumina
30 iSelect HD BeadChip. Genomic DNA and cDNA derived from six individuals were
31 sequenced on a 454 GS FLX system and generated a total of 1.2 million sequences, in which
32 SNPs were detected. As no reference genome exists for the house sparrow, we used the zebra
33 finch (*Taeniopygia guttata*) reference genome to determine the most likely position of each
34 SNP. The 10,000 SNPs on the SNP-chip were selected to be distributed evenly across 31
35 chromosomes, giving on average one SNP per 100,000 bp. The SNP-chip was screened across
36 1968 individual house sparrows from four island populations. Of the original 10,000 SNPs,
37 7413 were found to be variable, and 99% of these SNPs were successfully called in at least
38 93% of all individuals. We used the SNP-chip to demonstrate the ability of such genome-wide
39 marker data to detect population sub-division, and compared these results to similar analyses
40 using microsatellites. The SNP-chip will be used to map Quantitative Trait Loci (QTL) for
41 fitness-related phenotypic traits in natural populations.

42 **Introduction**

43 Single Nucleotide Polymorphisms (SNPs) have over the last decade been established as the
44 most frequently used genetic marker in genome mapping, population genomics and
45 Quantitative Trait Loci (QTL) mapping (Goddard & Hayes 2009; Mackay *et al.* 2009;
46 Allendorf *et al.* 2010; Slate *et al.* 2010). Until recently, SNPs have mainly seen use in studies
47 of the human genome and genetic model species, but - as the cost of next generation
48 sequencing decreases, SNPs are more readily used in studies of wild non-model species,
49 surpassing microsatellite markers as the major genetic marker technology (Luikart *et al.* 2003;
50 Kohn *et al.* 2006; Ekblom & Galindo; 2011). Compared to microsatellites, SNPs offer several
51 advantages including a higher frequency of occurrence in the genome in both coding and non-
52 coding regions (Collins *et al.* 1998; Brumfield *et al.* 2003), and a bi-allelic nature which
53 corresponds more closely with the models of evolution often applied in population genetics
54 (e.g. Hartl & Clark 1997). Additionally, SNPs offer logistical advantages such as a lower rate
55 of genotyping error and greater ease of automating large scale genotyping (Kennedy *et al.*
56 2003). The ability to efficiently type large numbers of SNP markers across numerous
57 individuals has enabled the study of the genetic architecture of quantitative traits by use of
58 QTL or gene mapping (Slate 2005; 2008; Goddard & Hayes 2009; Stapley *et al.* 2010).
59 Medium and high density genome-wide SNP panels are therefore powerful tools when
60 disentangling the relative roles of selection, genetic drift and gene flow in observed patterns
61 of genetic variation (Luikart *et al.* 2003; Nielsen *et al.* 2005; Stinchcombe & Hoekstra 2008).
62 Historically, QTL mapping has been performed mainly on laboratory populations and
63 outcrossed lines of cultured plants or domesticated animals under controlled conditions
64 (Duncan *et al.* 2007; Goddard & Hayes 2009). QTL mapping studies of natural populations
65 are more challenging than laboratory studies (Slate 2005; Slate *et al.* 2010), but are important

66 in order to establish the extent to which studies in model organisms can be extrapolated, and
67 to study adaptive genetic processes in systems where forces such as environmental
68 interactions, pleiotropy and epistasis have not been reduced or eliminated (e.g. Roff & Simons
69 1997; Kroymann & Mitchell-Olds 2005; Ellegren & Sheldon 2008). In order to perform QTL
70 mapping in a natural population, a large number of genetic markers and a detailed pedigree
71 are required. The number of genetic markers has previously been the major limiting factor for
72 QTL mapping in natural populations, but with the advent of applications derived from the
73 now frequently used next generation sequencing this is changing rapidly.

74 A method for *de-novo* development of a medium density SNP-chip for a non-model
75 species, the great tit (*Parus major*) has been described in Van Bers *et al.* (2012). *De-novo*
76 development of medium to high density SNP chips for non-model species is expected to
77 become more common over the next decade (Seeb *et al.* 2011), and best practice guides that
78 describe successful approaches will likely be important resources for future studies.
79 Additionally, comparisons of results obtained from SNP genotyping and those based on
80 microsatellites - which have a much higher mutation rate than SNPs (Foll & Gaggiotti 2008) -
81 are rare (Helyar *et al.* 2011). Here, we describe the development of a 10K Illumina iSelect HD
82 BeadChip for the house sparrow (*Passer domesticus*), a non-model species for which only a
83 small fraction of the genome has been sequenced but that is used to study many ecological,
84 evolutionary and physiological questions (e.g. Anderson 2006). The house sparrow is globally
85 distributed, either as a result of natural dispersal or from introductions mitigated by humans
86 (Anderson 2006). The ubiquitous nature of the house sparrow and multiple historical bottle-
87 neck events make it attractive for studies into evolution and adaptation (e.g. Summers-Smith
88 1988; Anderson, 2006). The sparrow breeds in proximity of human habitations, where they
89 nest in accessible cavities in houses, inside barns or within nest boxes. The house sparrow

90 lends itself as a good study species of adaptation because it has a short generation time
91 (approx. 2 years; Jensen *et al.* 2008), is easy to locate and capture, and can be measured and
92 non-lethally sampled without causing negative effects on its populations. Our study
93 complements that of Van Bers *et al.* (2012) by describing an alternative approach to a
94 successful 10K Illumina iSelect HD BeadChip *de-novo* development and large scale
95 population genotyping. Additionally, we investigated the performance of the SNP panel in
96 tests of population differentiation and compared the results to data derived from microsatellite
97 typing of the same populations.

98 The 10K Illumina iSelect HD BeadChip developed in this study will be used to
99 establish the first medium density marker map of the house sparrow genome, and will be an
100 important resource for QTL mapping and studies of the genetic architecture of complex traits,
101 genetic drift and adaptive evolution in both natural and manipulated populations.

102

103 **Materials and Methods**

104 Our data were collected from four islands within a natural house sparrow metapopulation in
105 northern Norway. The islands are separated by at least 20 km, are part of individual-based
106 long-term studies initiated in 1993 (Aldra [66°25'N, 13°04'E] and Hestmannøy [66°33'N,
107 12°50'E]) or 2001 (Leka [65°06'N, 11°38'E] and Vega [65°40'N, 11°55'E]) and used as
108 model systems to answer questions in evolutionary biology, ecology and conservation biology
109 (see e.g. Ringsby *et al.* 2002; Jensen *et al.* 2007; 2008; Pärn *et al.* 2009; Holand *et al.* 2011;
110 Billing *et al.* 2012; Pärn *et al.* 2012). Within this metapopulation system, less than 10% of
111 recruits disperse from their natal population (Pärn *et al.* 2009) and populations are
112 morphologically and genetically differentiated (Holand *et al.* 2011; Jensen *et al.* in review).
113 Genetic pedigrees going back > 8 generations are established for all island populations

114 represented in this study (Jensen *et al.* 2003; 2004; 2008; Billing *et al.* 2012; Rønning *et al.* in
115 prep.) and include approx. 6000 individuals.

116 *Collection of tissues and RNA extraction*

117 In February 2009, three relatively large populations (population size >160 adults) were
118 randomly sampled for one female and one male specimen. The populations included
119 Hestmannøy, Vega and Leka of the coast of central and northern Norway (Fig. 1). The house
120 sparrows were captured in mist nets, after which 47 µl of blood was immediately collected.
121 The sparrows were then quickly euthanized by cervical dislocation following guidelines
122 established and approved by the Norwegian Directorate for Nature Management, and within
123 10 minutes 80-150 mg of tissues from heart, liver, kidney, lung and brain were dissected.
124 Additionally, testis tissue (10 and 130 mg) was collected from two males, one from Leka and
125 the other from Vega. Blood and tissue samples were transferred to separate 2 ml DNase and
126 RNase free micro tubes (Nunc) containing RNAlater (QIAGEN, 1000 µl for blood samples
127 and 1600 µl for tissue samples) and immediately frozen for later extraction of RNA and DNA
128 in the lab. Total RNA from all tissues from one Hestmannøy male, testis samples from the
129 Leka and Vega Island males, and blood samples from all six individuals were extracted using
130 a RiboPure Blood Kit (Aambion Inc., Austin, TX, USA), with additional DNase treatment
131 following manufacturers' recommendations. For the remaining five individuals, total RNA
132 and DNA from all remaining tissues were extracted using a GeneMole automated nucleic acid
133 extraction system (Mole Genetics AS, Oslo, Norway) and the Total RNA Basic Kit (Mole
134 Genetics) without DNase treatment by Mole Genetics AS (Oslo, Norway). Extracted RNA
135 was stored in <2.5 mM Tris-HCL, pH 7.6.

136 *cDNA library generation and 454-sequencing*

137 For each individual, total RNA from all tissue samples and the blood sample was pooled prior to
138 cDNA library synthesis. The amount of extracted RNA for each tissue and blood was then
139 adjusted to give similar amounts of RNA from each, and 10 µg RNA in total from an individual.
140 One random-primed and normalized cDNA library was synthesized for each individual at
141 Vertis Biotechnologie AG (Friesing-Weihenstephan, Germany), following their in house
142 protocol for cDNA library generation. To prepare the cDNA for 454-Titanium sequencing,
143 normalized cDNA in the size range of 500 – 800 bp was eluted from preparative agarose gels.
144 The six synthesized cDNA libraries contained on average 17.2 ng/µl (range: 9-31 ng/µl) in 20
145 µl solutions. One run of 454-sequencing was carried out on a GS FLX Titanium (Roche,
146 Switzerland) at the Norwegian High-Throughput Sequencing Centre, University of Oslo,
147 Norway.

148 *SNP-chip development*

149 *Read mapping* Alignment of the reads was performed in CLC Genomics Workbench 4 (CLC
150 Bio, Aarhus, Denmark). All sequences from the six individuals were mapped together, thus
151 the contigs were allowed to contain sequences from more than one individual. Stringent
152 alignment parameters were used to ensure a high sequence quality from which to detect SNPs:
153 mismatch cost = 3; insertion cost = 3; deletion cost = 3; length fraction = 0.4; similarity =
154 0.99; ambiguity codes were used; non-specific matches were ignored and minimum contig
155 length set to 150 bp.

156 *SNP detection and selection* SNPs were detected in CLC Genomics using the following
157 parameters and Roche Phred scores: window length = 9; maximum number of gaps and
158 mismatches = 1; minimum average quality of surrounding bases = 20; minimum quality of
159 central base = 37; minimum coverage = 3; and minor allele frequency = 5%. Only non-

160 complex SNPs of the Infinium II design (one probe per SNP) with 60 bp of sequence on
161 either side of the SNP were chosen. The selected 121 bp sequences were searched against the
162 *Taeniopygia guttata* reference genomic sequences database build 1.1 on nucleotide BLAST
163 (<http://blast.ncbi.nlm.nih.gov>) as follows: optimised for discontinuous megablast; maximum
164 target sequences = 10; expected threshold = 10; word size = 11; match/mismatch scores = 1/-
165 1; gap cost = 2 for existence, 1 for extension and filtered for low complexity regions. Because
166 the 454-reads used to build contigs contained both genomic DNA and cDNA, the search
167 against *T. guttata* reference genomic sequences served to help us determine whether SNPs
168 from cDNA sequences were situated on intron/exon breaks and to inform us of the position of
169 the SNP on the *T. guttata* genome. The retained query sequences that contained intron/exon
170 breaks were then cropped at the intron/exon break and were not allowed to be less than 116 bp
171 long. Query sequences that had hits to more than one position on the *T. guttata* genome were
172 considered repeats and removed. Query sequences with BLAST e-values worse than $2.0E-15$
173 and those with more than six ambiguities were removed (i.e. the SNP itself and maximum 5
174 other ambiguities within the 116-121 bp sequence was allowed). In order to determine the
175 SNPs that were situated in known *T. guttata* mRNA regions, the remaining sequences were
176 put through BLAST against the *T. guttata* RNA reference genomic database
177 (<http://blast.ncbi.nlm.nih.gov>), with BLAST parameters as described above. A second BLAST
178 search against the more improved *T. guttata* reference genomic sequences database build
179 3.2.4.58 was carried out after development of the chip, in order to verify the positions on the
180 genome. The BLAST parameters were identical to previous searches. In total 19,852 SNPs
181 were sent to Illumina for processing by the Illumina Assay Design Tool to generate a score
182 file with a score and a failure code for each SNP that indicated the expected success for
183 designing an assay for the SNP. The SNPs that received failure codes equal to zero and scores

184 over 0.85 (N = 13,800) were retained. Of these, 9955 unique SNPs were selected to be
185 included on the 10K SNP chip based on their position on the *T. guttata* genome. The SNPs
186 were selected to be evenly distributed across the genome. The minimum distance between two
187 SNPs was set to 675 bp. For positions with more than one SNP to choose from, the one with
188 the highest read depth was selected. See Fig. 2 for an illustration of SNP distribution across
189 the *T. guttata* genome. As positive controls, 45 SNPs were randomly chosen to be typed twice
190 to test for genotyping errors.

191 The 10,000 SNPs included 10 SNPs located in candidate genes for beak morphology
192 and limb development: Calmodulin (Abzhanov *et al.* 2006; Schneider 2007), FGF8
193 (Abzhanov & Tabin 2004; Grant *et al.* 2006) and Frizzled (Brugmann *et al.* 2010). These
194 were specifically retained through the selection process despite not always conforming to the
195 criteria described above.

196 *Extraction of genomic DNA for sample screening*

197 The blood samples used for the SNP genotyping were collected from 1968 different
198 individuals on the islands Aldra (N = 406), Hestmannøy (N = 447), Leka (N = 512) and Vega
199 (N = 603) (see Fig. 1) during the years 1993-2010. The six individuals that were sequenced
200 (see above) were included among the 1968 samples. Of the 1968 individuals, 37 were
201 randomly chosen to be genotyped several times as positive controls. These 37 individuals
202 were typed two (N = 31), three (N = 5) or four (N = 1) times, thus the total number of DNA
203 samples screened was 2012. Additionally, four samples containing only ddH₂O were included
204 in the screening as negative controls. In total 2016 SNP chips were used in the genotyping.

205 Whole blood preserved in 96% ethanol, which had been stored for up to 10 years in
206 room temperature at the time of DNA extraction, was lysed in 60 µl Lairds buffer (Ausubel *et*
207 *al.* 1989), with 90 µg proteinase K (Sigma Aldrich, St Louis, MO), and incubated at 50°C for

208 3 hours. Genomic DNA was extracted from the lysate using the ReliaPrep Large Volume HT
209 gDNA Isolation System (Promega, Madison, WI), automated on a Biomek NXp robot
210 (Beckman Coulter, Miami, FL) and following the manufacturer's recommendations; the only
211 exception being elution of DNA in 25 mM Tris HCl (pH 8). The DNA concentrations were
212 measured using a Flurostar Omega scanner (MBG Labtech, Offenburg, Germany). Illumina
213 recommends sample concentrations of 40 – 60 ng/μl. Samples with concentrations above 60
214 ng/μl were normalized to a concentration of 50 ng/μl with 25 mM Tris HCl (pH 8). 760
215 samples had stock DNA concentrations below the recommended 40 ng/μl. For each sample, a
216 four μl aliquot of DNA containing a total of on average 160 ng (SD=56 ng) DNA was stored
217 at -20 °C until sample screening on the Illumina iSelect HD BeadChip. The SNP screening,
218 clustering and scoring of genotypes were carried out by the Genomics Core Facility,
219 Norwegian University of Science and Technology. The results were checked, filtered and
220 scored using GenomeStudio (Illumina, San Diego) following the guidelines provided by
221 Illumina (Illumina 2010).

222 *Descriptive statistics and analysis of genetic differentiation*

223 Descriptive statistics, tests regarding SNP design principles and SNP typing results were
224 carried out in the statistical software IBM SPSS 19.0 (SPSS Inc., 2010). The dataset used for
225 descriptive statistics was based on all loci for all genotyped samples except the four negative
226 controls. Quality control filtering of the dataset prior to analyses of genetic structure was done
227 in PLINK version 1.07 (<http://pngu.mgh.harvard.edu/purcell/plink/>; Purcell *et al.* 2007), using
228 a minor allele frequency (MAF) of 0.01 and maximum per-person missing (MIND) of 0.1.
229 Tests for Hardy-Weinberg equilibrium were done separately for each population and with a
230 significance level of 0.05. For the analyses of genetic structure and F_{ST} we used a reduced
231 dataset comprising individuals present on the four islands in 2002 (N = 173), and their

232 respective genotypes for SNPs that were found on autosomes and that passed the above
233 described quality control filtering ($N = 6736$). Individuals from 2002 were chosen in order to
234 avoid potential effects of experimental manipulation in two of the populations after 2002 and
235 to avoid inclusion of close relatives across multiple generations. First, a principal component
236 analysis (PCA) was performed using the R-package *adegenet* version 1.3 (Jombart 2008;
237 Jombart & Ahmed 2011) in the statistical software R version 2.14.2 (R Development Core
238 Team 2012). Second, genetic differentiation between populations was analysed with
239 STRUCTURE 2.3.3 (Pritchard *et al.* 2000) using no prior population information, the
240 admixture model, 10,000 burn-ins and 50,000 iterations, the number of populations from $K =$
241 1 to $K = 5$ and 10 separate runs for each K . The results were processed in STRUCTURE
242 HARVESTER (Earl & von Holdt 2012), which uses the Evanno method to determine the
243 most likely K (Evanno *et al.* 2005). Finally, we performed F_{ST} analyses in the R package
244 HIERFSTAT (Goudet 2005) which estimates F_{ST} with 95% confidence intervals (CI); if the
245 95% CI do not include zero the estimate is regarded as significantly different from zero at $P =$
246 0.05. The results from STRUCTURE and HIERFSTAT were compared to measures of
247 population divergence derived from microsatellite data in Jensen *et al.* (2013).

248

249 **Results**

250 *RNA extraction*

251 The RNA/DNA yield for samples extracted at Mole Genetics was 100 μ l eluate containing on
252 average 139.3 ng/ μ l (SD = 88.0 ng/ μ l) nucleic acids for liver, kidney, lung and brain tissue
253 samples, and on average 18.1 ng/ μ l (SD = 4.9 ng/ μ l) for heart tissue samples. For samples
254 extracted in-house the 100 μ l eluate contained on average 450.3 ng/ μ l (SD = 257.3 ng/ μ l)
255 total RNA for liver, kidney, lung and brain tissue samples, 27.5 ng/ μ l for the heart tissue

256 sample, 401.4 ng/μl (SD = 420.6 ng/μl) for the testis samples, and on average 52.3 ng/μl (SD
257 = 27.6 ng/μl) for the blood samples.

258 *SNP chip development*

259 *Sequencing* The 454-sequencing generated 1,160,122 reads of mean length 282 bp and
260 327,536,336 bp in total (Sequence read Archive accession numbers xxxxxxxx – xxxxxxxx). On
261 average, the number of reads from each of the six individuals was 191,242 (range: 145,889 –
262 287,470). The five individuals for which tissue samples were extracted at Mole Genetics had
263 length distribution of reads biased towards short reads (mean = 245 bp). This was due to
264 repeats in many of the sequences, which caused premature termination of the 454-sequencing.
265 The likely reason for such repeats was failure to treat the samples with DNase during RNA
266 extraction at Mole Genetics, with the consequence that genomic DNA was present in the
267 samples when cDNA libraries were synthesized. In contrast, for the individual that was
268 extracted in-house, the length distribution was as expected (i.e. a peak in the distribution of
269 read lengths at approx. 450 bp), indicating that sequences from this individual represented
270 RNA.

271 *Contig assembly and SNP detection* Stringent alignment of the reads produced 93,351 contigs
272 in which SNPs could be detected. The number of SNPs varied according to the stringency of
273 the search requirements. Using stringent search parameters described above but allowing
274 coverage down to 3 reads per site, we identified a total of 43,198 SNPs. Of these, 37,714
275 SNPs fulfilled the initial criteria of sequence length and SNP quality and were thus put
276 through BLAST. Of these, 13,800 SNPs satisfied the selection criteria (intron/exon breaks
277 were not allowed to extend more than a total of 5 bp into the 121 bp sequences, Infinium Type

278 II design, less than 5 ambiguities, a unique position on the *T. guttata* genome, and an Illumina
279 score of over 0.85) and were thus considered for inclusion on the chip.

280 *SNP-chip design characteristics and SNP call rates* Among the 10,000 SNPs selected for the
281 SNP-chip, the Illumina score was on average 0.9683, the mean number of ambiguities within
282 the 112-121 bp query sequence was 1.84 (range: 1 - 6, including the SNP itself), the mean
283 coverage was 6.37 (range: 3 - 67), and the mean count of the less frequent allele was 1.69
284 (range: 1 - 32) (Fig. 3).

285 Of the 10,000 SNPs on the chip, 8491 were successfully called (dbSNP accession
286 numbers xxxx – xxxx) and 7413 were variable. The Illumina design score, number of
287 ambiguities in the query sequence, coverage and the count of the least frequent allele was
288 similar for the 8491 SNPs that were successfully called and the 1509 SNPs that were not
289 called (Mann-Whitney tests: $P > 0.11$). On the other hand, the successfully called SNPs had
290 on average slightly longer query sequences than non-called SNPs (120.211 [SE = 0.014] vs.
291 120.123 (SE = 0.036), Mann-Whitney test: $P = 0.015$) and worse e-values (2.0E-18 [SE =
292 7.7E-19] vs. 5.5E-19 [SE = 2.2E-19], Mann-Whitney test: $P < 0.001$). These results suggest
293 that few of the SNP-chip design criteria affected the probability of a SNP to be successfully
294 called (given that the SNP had fulfilled our detection and selection criteria).

295 Compared to the 1078 called and non-variable SNPs, the 7413 variable SNPs had on
296 average slightly lower read depth (6.34 [SE = 0.05] vs. 6.67 [SE = 0.13], Mann-Whitney tests:
297 $P = 0.002$) but higher count of the least frequent allele in the sequence data (1.78 [SE = 0.02]
298 vs. 1.14 (SE = 0.01], Table 1, Mann-Whitney tests: $P < 0.001$). This suggests that a SNP was
299 more likely to be a true SNP (i.e. a polymorphic base) if the number of sequences in which its
300 rarest allele was observed was two rather than one.

301 *Assessment of SNP call rate and genotyping error in genotyped samples*

302 Of the 2016 samples that were genotyped (including four negative controls and the six
303 sequenced individuals), SNPs were called for 1999 samples. The 13 samples besides the
304 negative controls that were not called had on average lower DNA concentration (mean: 17.79
305 ng/ μ l, SD = 18.20 ng/ μ l, range: 3.99 - 50 ng/ μ l) than the 1999 samples that were called
306 (mean: 40.05 ng/ μ l, SD = 13.83 ng/ μ l, range: 0.83 - 60.00 ng/ μ l) (Mann-Whitney test: $P <$
307 0.001). For samples with a DNA concentration lower than 20 ng/ μ l, the probability of
308 successful calling of SNPs decreased; no SNPs were called for 10 out of 284 samples (3.5%)
309 with a DNA concentration lower than 20 ng/ μ l. In contrast, no SNPs were called for only 3
310 out of 1728 samples (0.2%) with a DNA concentration of 20 ng/ μ l or more. From the 37
311 individuals that were typed twice or more, we estimated a SNP typing error rate of 0.0005%.
312 Additionally, 45 SNPs were included twice on the chip: 32 SNPs returned the same genotype
313 for all pairs (typed in between 1995 and 1999 samples), whilst 10 SNPs returned a genotype
314 for one of the assays but failed for the other. For three SNP pairs no individuals were called.
315 Thus there were no conflicts among the 45 SNPs that were run twice, although some assays
316 did not return a genotype.

317 The mean number of SNPs called for a given sample was 8457 (range: 8169 - 8476).
318 The 7413 variable SNPs were on average typed in 1988 samples (Table 1). The mean minor
319 allele frequency of variable SNPs was 0.2380, and ranged from 0.00025 (i.e. one sample was
320 heterozygous at the SNP and the rest homozygous for the common base) to 0.5 (Table 1). Of
321 the variable SNPs, only 21.4% had a minor allele frequency below 0.1, suggesting that most
322 of the variable SNPs will be valuable in further analyses.

323 *Genetic differentiation of sub-populations*

324 Principal component analysis using SNP data from individuals sampled in 2002 indicated
325 three distinct clusters; the individuals from Aldra constituted one group, as did the individuals
326 from Leka, whilst Vega and Hestmannøy clustered together (but with incomplete overlap)
327 (Fig. 4). In concordance with these results, F_{ST} analyses indicated that Aldra was the most
328 differentiated population, with pair-wise values approximately twice as high as for the other
329 islands (Table 2). This pattern was consistent also in the STRUCTURE analysis, which under
330 the most likely scenario identified two clusters; again with Aldra as a distinct population and
331 the remaining islands clustering together (Fig. 5). Under the less likely scenario of three
332 clusters, the pattern was the same as for the principal component analysis, with Aldra and
333 Leka being separate groups and Hestmannøy and Vega grouping together. Results based on
334 SNP data corresponded closely with results derived from microsatellites, which showed
335 similar levels of F_{ST} between pairs of populations (Table 2) and that - based on analyses in
336 STRUCTURE - the four islands fell into two different clusters; again with Aldra in one
337 cluster and the three remaining islands in a second group (Jensen *et al.* 2013).

338

339 **Discussion**

340 We have described the development of a 10K Illumina iSelect HD BeadChip for the house
341 sparrow, and have assessed the performance of the genome-wide SNP data to detect and
342 quantify genetic sub-division of a meta-population in central and northern Norway. The
343 available genomic resources derived from our study include an additional 30,000 putative
344 house sparrow SNPs in 93,351 contigs. With a mean length of 475 bp, these contigs cover
345 approximately 44.5 million bp of the house sparrow genome, and feature both coding and
346 non-coding regions. Assuming that the house sparrow genome has the same size as the zebra
347 finch genome (1.2 Gbp; Warren *et al.* 2010), these sequences cover about 3.7% of the total

348 genome. This resource will be important during future development of SNP-chips for further
349 investigation of specific QTL regions, or for primer design in a candidate gene approach
350 where desired target genes are partly or wholly included in the contigs. Additionally,
351 considering the cross-population success in great tits described in Van Bers *et al.* (2012), we
352 predict that the house sparrow SNP chip will provide reliable results if applied to house
353 sparrow populations outside of Norway, as indicated by microsatellite genetic structure for
354 house sparrow populations distributed globally (Schrey *et al.* 2011). Moreover, successful
355 cross species applications of a 50K *Ovis aries* SNP chip have been described for bighorn and
356 thinhorn sheep, with call rates of 95 and 90% respectively (Miller *et al.* 2011). It is therefore
357 likely that the house sparrow SNP chip could be applied also to other species or sub-species
358 of the European *Passer* genus, which has been found to have multiple hybrid zones and pair-
359 wise F_{ST} estimates that are comparable with those found on house sparrows along the
360 Norwegian coast (Hermansen *et al.* 2011; Jensen *et al.* 2013).

361 The approach we have described for SNP-chip development proved highly successful.
362 The overall SNP call rate was 85% and approximately 75% of the SNPs were variable and
363 informative for population differentiation and/or marker map development. Of our 2012
364 samples, 99.35% were called for at least 8169 SNPs. The error rate estimated from duplicate
365 samples and SNPs was very low (0.0005%). For comparison, the call rate expected for
366 application of commercially available SNP-chips developed for humans is approximately
367 98% (International HapMap Consortium 2010). Other studies on model species have for
368 example reported call rates of 93%, with 89% of the SNPs being polymorphic for a 60K
369 chicken SNP chip (Groenen *et al.* 2011). For non-model species the success rate is generally
370 somewhat lower: when genotyping wild and farmed Atlantic salmon (*Salmo salar*) on a 7K
371 SNP-chip, Karlsson *et al.* (2011) reported that 65% of the SNPs were called and informative,

372 whilst Van Bers *et al.* (2012) developed a 10K SNP chip for two populations of the great tit
373 and obtained a call rate of 83%, with about 72% of the SNPs being polymorphic. Our custom
374 made 10K SNP-chip thus has a genotyping success rate that is somewhat lower than
375 commercial SNP-chips for humans but has the same proportion of called and variable SNPs
376 as other studies of non-model species.

377 The inclusion of SNPs located in genomic DNA was un-intentional and caused by
378 failure to treat tissue samples collected from five of the six sequenced individuals with
379 DNase. The original idea was to sequence cDNA to obtain a high read depth in which to
380 detect SNPs, with relatively low sequencing costs, and avoid choosing SNPs with very low
381 minor allele count. In hindsight, it seems that sequencing a mix of cDNA and gDNA did not
382 reduce the quality of our SNP-chip compared to similar studies, despite the fact that we had to
383 choose some SNPs that were detected in regions with only 3x read depth.

384 A large number of reads at the SNP sites is advantageous, as it allows for greater
385 confidence that the SNP is not an artifact of sequencing error. For their 60K chicken SNP-
386 chip, Groenen *et al.* (2011) used read depth of 12x or more and Van Bers *et al.* (2012)
387 reported a read depth at the SNP site of > 8. The median read depth of the SNPs included on
388 the house sparrow chip was 5 (range 3 – 67), whilst 2569 of the SNPs were situated in regions
389 with a read depth of 3 (see Fig. 3). Of these, 60% (N = 1552) were both successfully called,
390 variable both within and across our sampled populations and with a minor allele frequency \geq
391 0.01. Comparably, this is a lower success rate than for SNPs with a higher read depth, but
392 indicates that SNPs detected in low-coverage regions have the potential to be highly useful.

393 Approximately one third of our samples had a lower than recommended (< 40 ng/ μ l)
394 DNA concentration, however genotyping success was only marginally affected: there was a
395 15-20% lower call rate for samples with concentrations below 10 ng/ μ l, and a 1% lower call

396 rate for samples with concentrations ranging from 10-20 ng/ μ l. It therefore appears that the
397 Illumina iSelect HD BeadChip requires less sample material than some sequencing based
398 SNP genotyping techniques (Miller *et al.* 2007, Baird *et al.* 2008). The approach we have
399 developed may therefore be a useful resource for molecular ecology studies on organisms for
400 which only small amounts of DNA can be acquired.

401 In all analyses of genetic population structure, the island of Aldra was the most
402 divergent, a pattern which is probably explained by the unique history of this population.
403 Aldra is known to have been colonized by four founders in 1998 and has since received few
404 immigrants (Billing *et al.* 2012). The effective population size on this island has been
405 estimated to range from approximately 10 to 30 (Engen *et al.* 2007; Baalsrud *et al.* in review).
406 Accordingly, the level of inbreeding in this island population is significant (Jensen *et al.*
407 2007; Billing *et al.* 2012), and it is likely that an initial founder effect and subsequent genetic
408 drift can explain the strong genetic divergence of Aldra from the other islands.

409 The results from STRUCTURE and F_{ST} analyses indicated that our panel of approx.
410 7000 variable SNPs and a panel of 14 highly variable microsattellites produced very similar
411 results. A denser panel of SNPs is required in order to achieve the same power as
412 microsattellites (Evans & Cardon 2004). For instance Hess *et al.* (2011) found that
413 microsattellites performed better than SNPs when 13 highly variable microsatellite loci were
414 compared to a SNP-panel of 92 loci for fine-scale population identification. Microsattellites
415 are in the process of being replaced by SNPs in a variety of molecular applications, including
416 those within molecular ecology (Allendorf *et al.* 2010; Ekblom & Galindo 2011; Hess *et al.*
417 2011). However, the costs for development of SNP panels with a power comparable to
418 microsattellites are still significant and may not be feasible for projects focusing on the
419 ecology of non-model species. Provided microsattellites in a panel have low error rates, are

420 independent and distributed across the genome, there do not seem to be strong reasons to
421 abandon the use of microsatellites for the purpose of investigating population differentiation
422 and assignment analysis. However, in many other cases, such as in QTL mapping, studies of
423 adaptive evolution and effects of genetic drift, a greater number of markers is required and
424 microsatellites fall short compared to high density SNP panels.

425 In conclusion, we have described an easy and cost effective protocol for successful
426 generation and population scale screening of a 10K medium density SNP chip in a non-model
427 species. We have shown that despite features such as lower than desired read depth (3, for
428 approximately 25% of the SNPs), lower than recommended DNA concentration for some
429 samples, and a large proportion of SNPs situated in genomic DNA, the result was very
430 successful and comparable with other medium density SNP-chip population screens on non-
431 model species. Lastly, our study indicated that for the purpose of population assignment and
432 differentiation, high density SNP data produce results that are comparable with those derived
433 from high quality microsatellite data. In the near future, the SNP data will be used to develop
434 the first marker map for house sparrows, and subsequent analysis into the genetic architecture
435 of quantitative traits in wild populations of house sparrows.

436

437 **Acknowledgements**

438 We are grateful to the hospitable and friendly inhabitants in the field area who made the study
439 possible. We also thank fieldworkers and laboratory technicians for assistance. We thank the
440 Research Council of Norway (Project no: 191847, Strategic University Program (SUP) in
441 Conservation Biology) and the functional genomics programme at the Norwegian University
442 of Science and Technology (to IJH and AMB) for funding. The SNP microarray service was
443 provided by the Genomics Core Facility, Norwegian University of Science and Technology - a

444 national technology platform supported by the functional genomics program (FUGE) of the
445 Research Council of Norway. Permit to collect tissue samples from 6 individual house
446 sparrows was given from the Norwegian Directorate for Nature Management. The research
447 was carried out in accordance with permits from the Norwegian Directorate for Nature
448 Management and the Bird Ringing Centre at Stavanger Museum, Norway.

449

450 **Data accessibility**

451 All sequence reads have been submitted to the Sequence read Archive (SRA). Accession
452 numbers are xxxxxxxx – xxxxxxxx. All genotyped SNPs with the 116-121 bp flanking
453 sequences have been submitted to dbSNP. Accession numbers are xxxxxxxx – xxxxxxxx.

454

455 **References**

- 456 Abzhanov A, Kuo WP, Hartmann C, *et al.* (2006) The calmodulin pathway and evolution of
457 elongated beak morphology in Darwin's finches. *Nature* **442**, 563-567.
- 458 Abzhanov A, Tabin CJ (2004) Shh and Fgf8 act synergistically to drive cartilage outgrowth
459 during cranial development. *Developmental Biology* **273**, 134-148.
- 460 Allendorf FW, Hohenlohe PA, Luikart G (2010) Genomics and the future of conservation
461 genetics. *Nature Reviews Genetics* **11**, 697-709.
- 462 Anderson TR (2006) Biology of the ubiquitous house sparrow: from genes to populations.
463 Oxford University Press, New York.
- 464 Ausubel FM, Brent R, Kingston RE, Moore DD, Seidman JG, Smith JA, Struhl K, eds. 1989.
465 Current Protocols in Molecular Biology. John Wiley & Sons, New York.
- 466 Baird N, Etter P, Atwood T, *et al.* (2008) Rapid SNP Discovery and Genetic Mapping Using
467 Sequenced RAD Markers. *PloS one* **3**.
- 468 Billing AM, Lee AM, Skjelseth S, *et al.* (2012) Evidence of inbreeding depression but not
469 inbreeding avoidance in a natural house sparrow population. *Molecular Ecology* **21**,
470 1487-1499.
- 471 Brugmann SA, Powder KE, Young NM, *et al.* (2010) Comparative gene expression analysis
472 of avian embryonic facial structures reveals new candidates for human craniofacial
473 disorders. *Human molecular genetics* **19**, 920-930.
- 474 Brumfield RT, Beerli P, Nickerson DA, Edwards SV (2003) The utility of single nucleotide
475 polymorphisms in inferences of population history. *Trends in Ecology and Evolution*
476 **18**, 249-256.
- 477 Collins FS, Brooks LD, Chakravarti A (1998) A DNA Polymorphism Discovery Resource for
478 Research on Human Genetic Variation. *Genome Research* **8**, 1229-1231.

- 479 Duncan EJ, Dodds KG, Henry HM, Thompson MP, Phua SH (2007) Cloning, mapping and
480 association studies of the ovine ABCG2 gene with facial eczema disease in sheep.
481 *Animal Genetics* **38**, 126-131.
- 482 Earl D, von Holdt B (2012) STRUCTURE HARVESTER: a website and program for
483 visualizing STRUCTURE output and implementing the Evanno method. *Conservation*
484 *Genetics Resources* **4**, 359-361.
- 485 Ekblom R, Galindo J (2011). Applications of next generation sequencing in molecular
486 ecology of non-model organisms. *Heredity* **107**, 1-15.
- 487 Ellegren H, Sheldon BC (2008). Genetic basis of fitness differences in natural populations.
488 *Nature* **452**, 169-175.
- 489 Engen S, Ringsby TH, Sæther B-E, *et al.* (2007) Effective size of fluctuating populations with
490 two sexes and overlapping generations. *Evolution* **61**, 1873-1885.
- 491 Evanno G, Regnaut S, Goudet J (2005) Detecting the number of clusters of individuals using
492 the software structure: a simulation study. *Molecular Ecology* **14**, 2611-2620.
- 493 Evans DM, Cardon LR (2004) Guidelines for Genotyping in Genomewide Linkage Studies:
494 Single-Nucleotide–Polymorphism Maps Versus Microsatellite Maps. *The American*
495 *Journal of Human Genetics* **75**, 687-692.
- 496 Foll M, Gaggiotti O (2008) A genome-scan method to identify selected loci appropriate for
497 both dominant and codominant markers: a Bayesian perspective. *Genetics* **180**, 977-
498 993.
- 499 Goddard ME, Hayes BJ (2009) Mapping genes for complex traits in domestic animals and
500 their use in breeding programmes. *Nature Reviews Genetics* **10**, 381-391.
- 501 Goudet J 2005 HIERFSTAT, a package for R to compute and test hierarchical F-statistics.
502 *Molecular Ecology Notes* **5**, 184–186.

- 503 Grant PR, Grant BR, Abzhanov A (2006) A developing paradigm for the development of bird
504 beaks. *Biological Journal of the Linnean Society* **88**, 17-22.
- 505 Groenen M, Megens H-J, Zare Y, *et al.* (2011) The development and characterization of a
506 60K SNP chip for chicken. *BMC Genomics* **12**, 274.
- 507 Hartl DL, Clark AG (1997) Principles of population genetics. 3rd ed. Sinauer Associates,
508 Sunderland.
- 509 Helyar SJ, Hemmer-Hansen J, Bekkevold D, *et al.* (2011) Application of SNPs for population
510 genetics of nonmodel organisms: new opportunities and challenges. *Molecular*
511 *Ecology Resources* **11**, 123-136.
- 512 Hermansen JS, Sæther SA, Elgvin TO, *et al.* (2011) Hybrid speciation in sparrows I:
513 phenotypic intermediacy, genetic admixture and barriers to gene flow. *Molecular*
514 *Ecology* **20**, 3812-3822.
- 515 Hess JE, Matala AP, Narum SR (2011) Comparison of SNPs and microsatellites for fine-scale
516 application of genetic stock identification of Chinook salmon in the Columbia River
517 Basin. *Molecular Ecology Resources* **11**, 137-149.
- 518 Holand AM, Jensen H, Tufto J, Moe R (2011). Does genetic drift or selection explain
519 geographic differentiation of morphological characters in house sparrows? *Genetics*
520 *Research* **93**, 367-379.
- 521 Illumina, Inc. (2010) Infinium genotyping data analysis: a guide for analyzing infinium
522 genotyping data using the Illumina GenomeStudio genotyping module. Illumina, Inc,
523 San Diego. Available at:
524 http://www.illumina.com/Documents/products/technotes/technote_infinium_genotypin
525 [g_data_analysis.pdf](http://www.illumina.com/Documents/products/technotes/technote_infinium_genotyping_data_analysis.pdf)

- 526 International HapMap Consortium (2010) Integrating common and rare genetic variation in
527 diverse human populations. *Nature* **467**, 52-58.
- 528 Jensen H, Sæther B-E, Ringsby TH, *et al.* (2003) Sexual variation in heritability and genetic
529 correlations of morphological traits in house sparrow (*Passer domesticus*). *Journal of*
530 *Evolutionary Biology* **16**, 1296-1307.
- 531 Jensen H, Sæther B-E, Ringsby TH, *et al.* (2004) Lifetime reproductive success in relation to
532 morphology in the house sparrow *Passer domesticus*. *Journal of Animal Ecology* **73**,
533 599-611.
- 534 Jensen H, Bremset EM, Ringsby TH, Sæther B-E (2007) Multilocus heterozygosity and
535 inbreeding depression in an insular house sparrow metapopulation. *Molecular Ecology*
536 **16**, 4066-4078.
- 537 Jensen H, Steinsland I, Ringsby TH, Sæther B-E (2008) Evolutionary dynamics of a sexual
538 ornament in the house sparrow (*Passer domesticus*): the role of indirect selection
539 within and between sexes. *Evolution* **62**, 1275-1293.
- 540 Jensen H, Moe R, Hagen IJ, Holand AM, Kekkonen J, Tufto J, Sæther B-E (2013) Genetic
541 variation and structure of house sparrow populations: is there an island effect?
542 *Molecular Ecology* **xx**, xx-xx.
- 543 Jombart T (2008) adegenet: a R package for the multivariate analysis of genetic markers.
544 *Bioinformatics (Oxford, England)* **24**, 1403-1405.
- 545 Jombart T, Ahmed I (2011) adegenet 1.3-1: new tools for the analysis of genome-wide SNP
546 data. *Bioinformatics* **27**, 3070-3071.
- 547 Karlsson S, Moen T, Lien S, Glover KA, Hindar K (2011) Generic genetic differences
548 between farmed and wild Atlantic salmon identified from a 7K SNP-chip. *Molecular*
549 *Ecology Resources* **11**, 247-253.

- 550 Kennedy GC, Matsuzaki H, Dong S, *et al.* (2003) Large-scale genotyping of complex DNA.
551 *Nature Biotechnology* **21**, 1233-1237.
- 552 Kohn MH, Murphy WJ, Ostrander EA, Wayne RK (2006). Genomics and conservation
553 genetics. *Trends in Ecology and Evolution* **21**, 629-637.
- 554 Kroymann J, Mitchell-Olds T (2005) Epistasis and balanced polymorphism influencing
555 complex trait variation. *Nature* **435**, 95-98.
- 556 Luikart G, England PR, Tallmon D, Jordan S, Taberlet P (2003) The power and promise of
557 population genomics: from genotyping to genome typing. *Nature Reviews Genetics* **4**,
558 981-994.
- 559 Mackay T FC, Stone EA, Ayroles JF (2009) The genetics of quantitative traits: challenges and
560 prospects. *Nature Reviews Genetics* **10**, 565-577.
- 561 Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA (2006) Rapid and cost-effective
562 polymorphism identification and genotyping using restriction site associated DNA
563 (RAD) markers. *Genome Research* **17**.
- 564 Miller JM, Poissant J, Kijas JW, Coltman DW, the International Sheep Genomics C (2011) A
565 genome-wide set of SNPs detects population substructure and long range linkage
566 disequilibrium in wild sheep. *Molecular Ecology Resources* **11**, 314-322.
- 567 Nielsen R, Williamson S, Kim Y, *et al.* (2005) Genomic scans for selective sweeps using SNP
568 data. *Genome Research* **15**, 1566-1575.
- 569 Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using
570 multilocus genotype data. *Genetics* **155**, 945-959.
- 571 Purcell S, Neale B, Todd-Brown K, *et al.* (2007) PLINK: A Tool Set for Whole-Genome
572 Association and Population-Based Linkage Analyses. *The American Society of Human*
573 *Genetics* **81**, 559-575.

- 574 Pärn H, Jensen H, Ringsby TH, Sæther B-E (2009) Sex-specific fitness correlates of dispersal
575 in a house sparrow metapopulation. *Journal of Animal Ecology* **78**, 1216-1225.
- 576 Pärn H, Ringsby TH, Jensen H, Sæther B-E (2012) Spatial heterogeneity in the effects of
577 climate and density-dependence on dispersal in a house sparrow metapopulation.
578 *Proceedings of the Royal Society B: Biological Sciences* **279**, 144-152.
- 579 Ringsby TH, Sæther B-E, Tufto J, Jensen H, Solberg EJ (2002) Asynchronous spatiotemporal
580 demography of a house sparrow metapopulation in a correlated environment. *Ecology*
581 **83**, 561-569.
- 582 R Development Core Team (2012) R: A language and environment for statistical computing.
583 R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL
584 <http://www.R-project.org/>.
- 585 Roff DA, Simons AM (1997) The quantitative genetics of wing dimorphism under laboratory
586 and field conditions in the cricket *Gryllus pennsylvanicus*. *Heredity* **78**, 235-240.
- 587 Schneider RA (2007) How to tweak a beak: molecular techniques for studying the evolution
588 of size and shape in Darwin's finches and other birds. *BioEssays* **29**, 1-6.
- 589 Schrey AW, Grispo M, Awad M, *et al.* (2011) Broad-scale latitudinal patterns of genetic
590 diversity among native European and introduced house sparrow (*Passer domesticus*)
591 populations. *Molecular Ecology* **20**, 1133-1143.
- 592 Seeb JE, Carvalho G, Hauser L, *et al.* (2011) Single-nucleotide polymorphism (SNP)
593 discovery and applications of SNP genotyping in nonmodel organisms. *Molecular*
594 *Ecology Resources* **11**, 1-8.
- 595 Slate J (2005) INVITED REVIEW: Quantitative trait locus mapping in natural populations:
596 progress, caveats and future directions. *Molecular Ecology* **14**, 363-379.

- 597 Slate J (2008) Robustness of linkage maps in natural populations: a simulation study.
598 *Proceedings. Biological sciences / The Royal Society* **275**, 695-702.
- 599 Slate J, Santure AW, Feulner PGD, *et al.* (2010) Genome mapping in intensively studied
600 vertebrate populations. *Trends in Genetics* **26**, 275-284.
- 601 Stapley J, Reger J, Feulner PGD, *et al.* (2010) Adaptation genomics: the next generation.
602 *Trends in Ecology & Evolution* **25**, 705-712.
- 603 Stinchcombe JR, Hoekstra HE (2008) Combining population genomics and quantitative
604 genetics: finding the genes underlying ecologically important traits. *Heredity* **100**,
605 158-170.
- 606 Summers-Smith JD (1988) The sparrows: a study of the genus *Passer*. T & AD Poyser Ltd,
607 Calton.
- 608 Van Bers NEM, Santure AW, Van Oers K, *et al.* (2012) The design and cross-population
609 application of a genome-wide SNP chip for the great tit *Parus major*. *Molecular*
610 *Ecology Resources* **12**, 753-779.
- 611 Warren WC, Clayton DF, Ellegren H, *et al.* (2010) The genome of a songbird. *Nature*, **464**,
612 757-762.
- 613

614 **Tables**

615

616 **Table 1:** Characteristics of SNPs on the 10K Illumina custom SNP-chip for house sparrows when typed in 2012 samples.

617

	Number of SNPs	Minor allele frequency	
		Mean	Median (min - max)
All SNPs	10000	0.197	0.1803 (0 - 0.5)
Non-called SNPs	1509	-	-
Called SNPs	8491	0.208	0.1958 (0 - 0.5)
Non-variable SNPs	1078	0	0 (0 - 0)
Variable SNPs	7413	0.238	0.2286 (0.00025 - 0.5)

618 **Table 2:** Pair-wise genetic distances between four island populations of house sparrows in
 619 Norway. Values below the diagonal are F_{ST} values based on 14 microsatellite loci (Jensen *et*
 620 *al.* 2013); values above diagonal are F_{ST} values derived from 6736 SNPs. F_{ST} values and
 621 their 95% confidence limits (in parentheses) were calculated using the R-package
 622 HIERFSTAT (Goudet 2005).
 623

	Hestmannøy	Aldra	Vega	Leka
Hestmannøy	-	0.073 (0.071-0.076)	0.024 (0.022-0.025)	0.029 (0.028-0.031)
Aldra	0.062 (0.044-0.083)	-	0.078 (0.076-0.081)	0.082 (0.079-0.085)
Vega	0.024 (0.016-0.034)	0.069 (0.041-0.101)	-	0.028 (0.027-0.029)
Leka	0.023 (0.013-0.036)	0.074 (0.053-0.096)	0.025 (0.016-0.034)	-

624
 625

626 **Figure legends**

627 **Figure 1:** Map showing the study area on the Norwegian coast. The four study islands are
628 shown in black.

629

630 **Figure 2:** Distribution of SNPs present on the 10K house sparrow SNP chip when mapped
631 onto the *T. guttata* genome (Tgu chromosomes 1 – 28 and Z). The number of SNPs within
632 500,000 bp intervals is shown.

633

634 **Figure 3:** Plots showing sequencing information for the 10,000 SNPs on the house sparrow
635 SNP chip. Light grey bars (left axis) show the number of SNPs on the chip selected in
636 genomic regions with a given sequence coverage. Dots (right axis) show mean count (\pm 1SD)
637 of the minor allele for each level of sequence coverage at the SNP location.

638

639 **Figure 4:** Plots of principle component analysis (PCA) of genetic variation between 173 adult
640 individuals present in 2002 in four island house sparrow populations off the coast of Norway.
641 The colours of each individual represent island population: Aldra = red, Hestmannøy = black,
642 Leka = green and Vega = blue. The PCAs are based on 6736 autosomal SNPs.

643

644 **Figure 5:** Structure barplots based on 6736 autosomal SNPs from 173 adult individuals
645 present in 2002 in four island house sparrow populations off the coast of Norway. The upper
646 panel shows results for $K = 2$, whilst the lower panel shows results for $K = 3$.

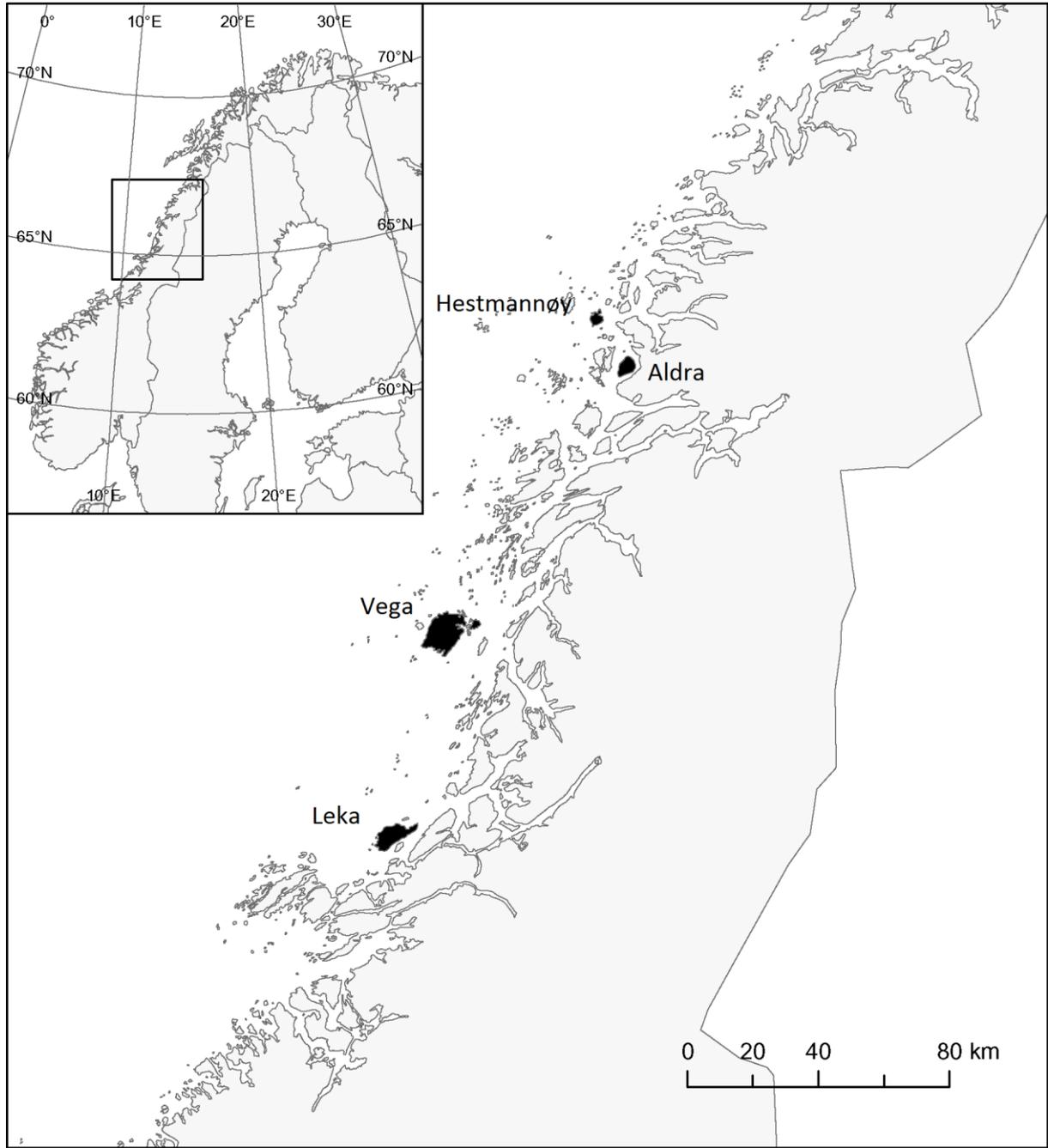
647

648

649 **Figures**

650

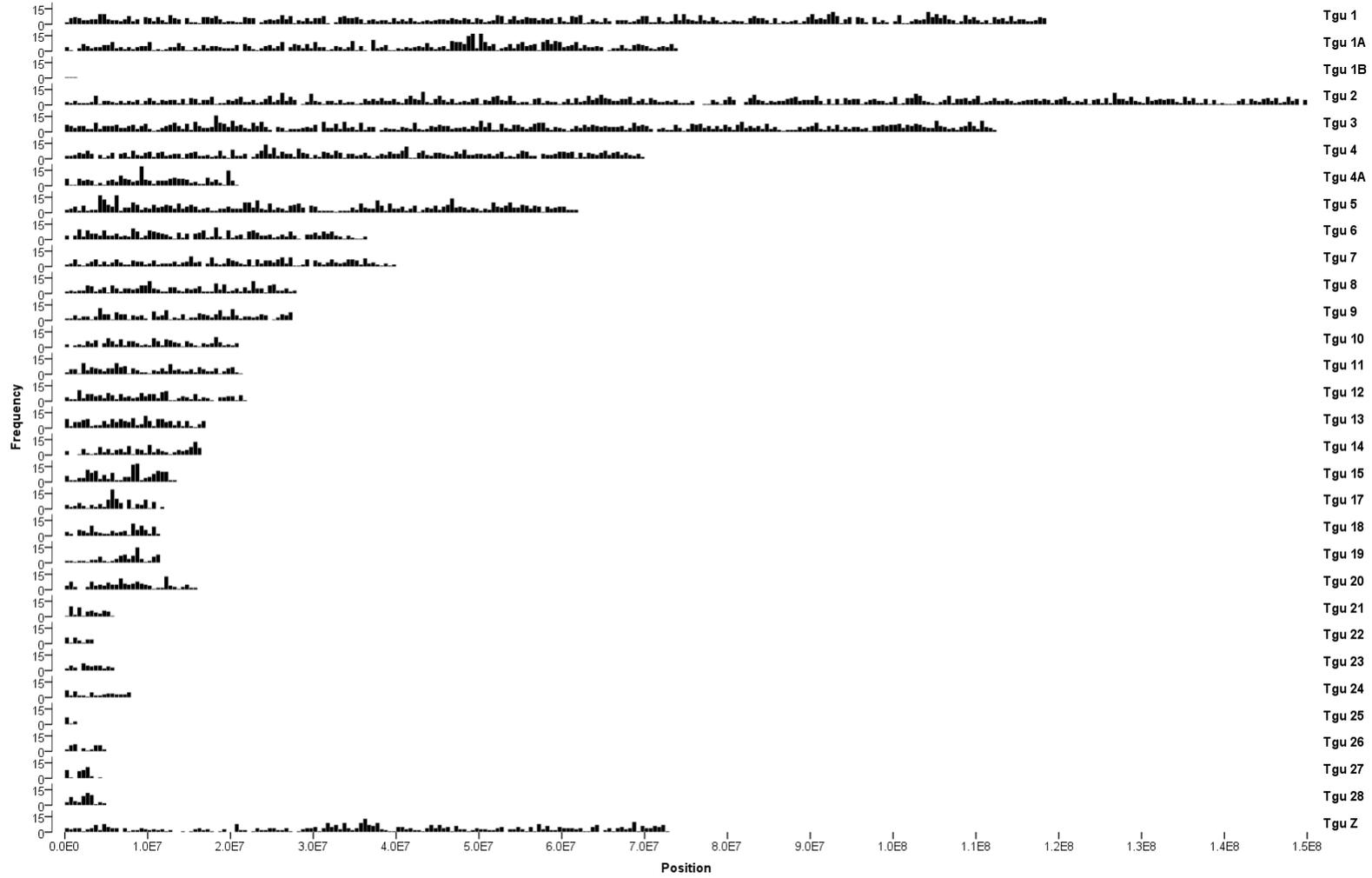
651 **Figure 1**



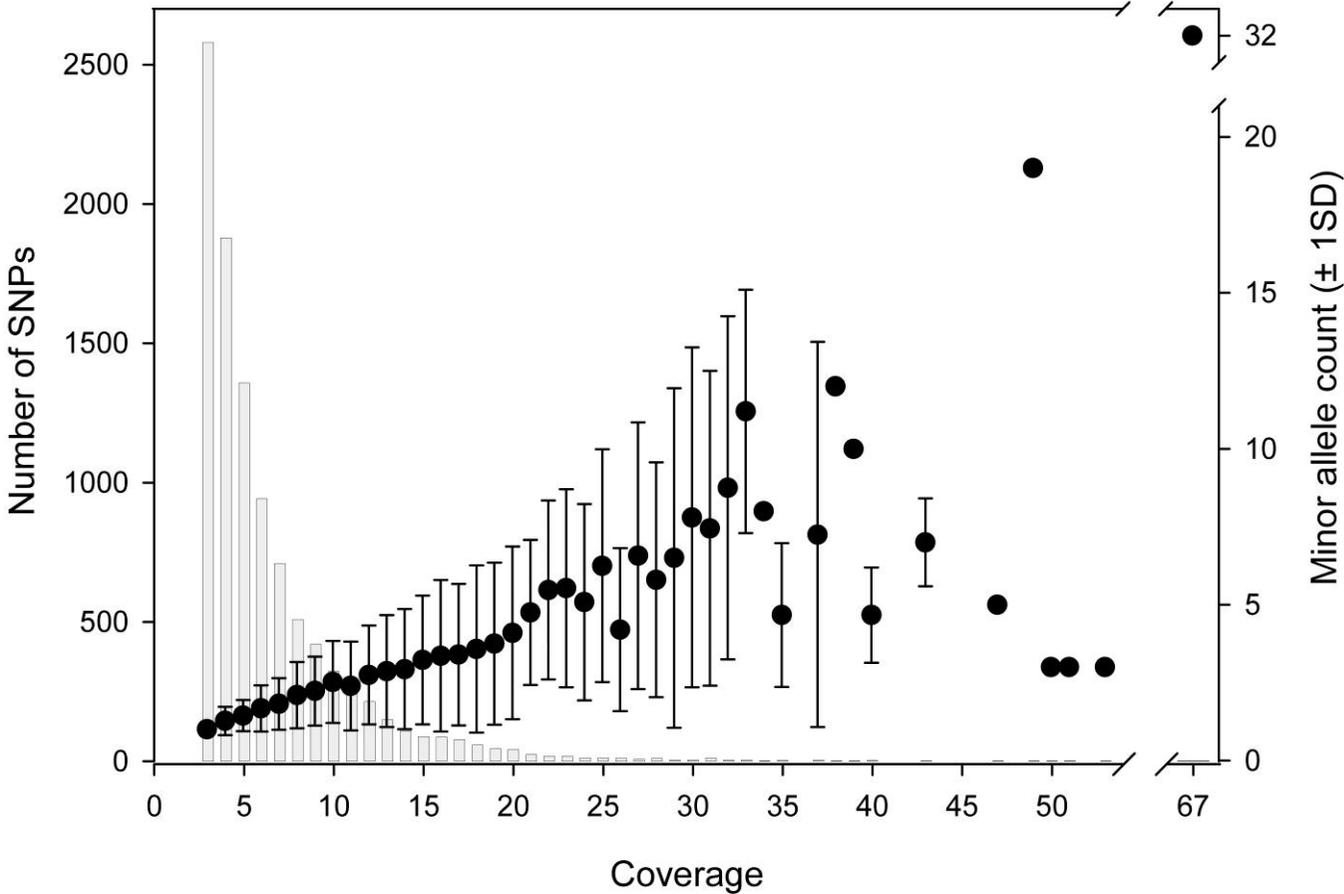
652

653 **Figure 2**

654

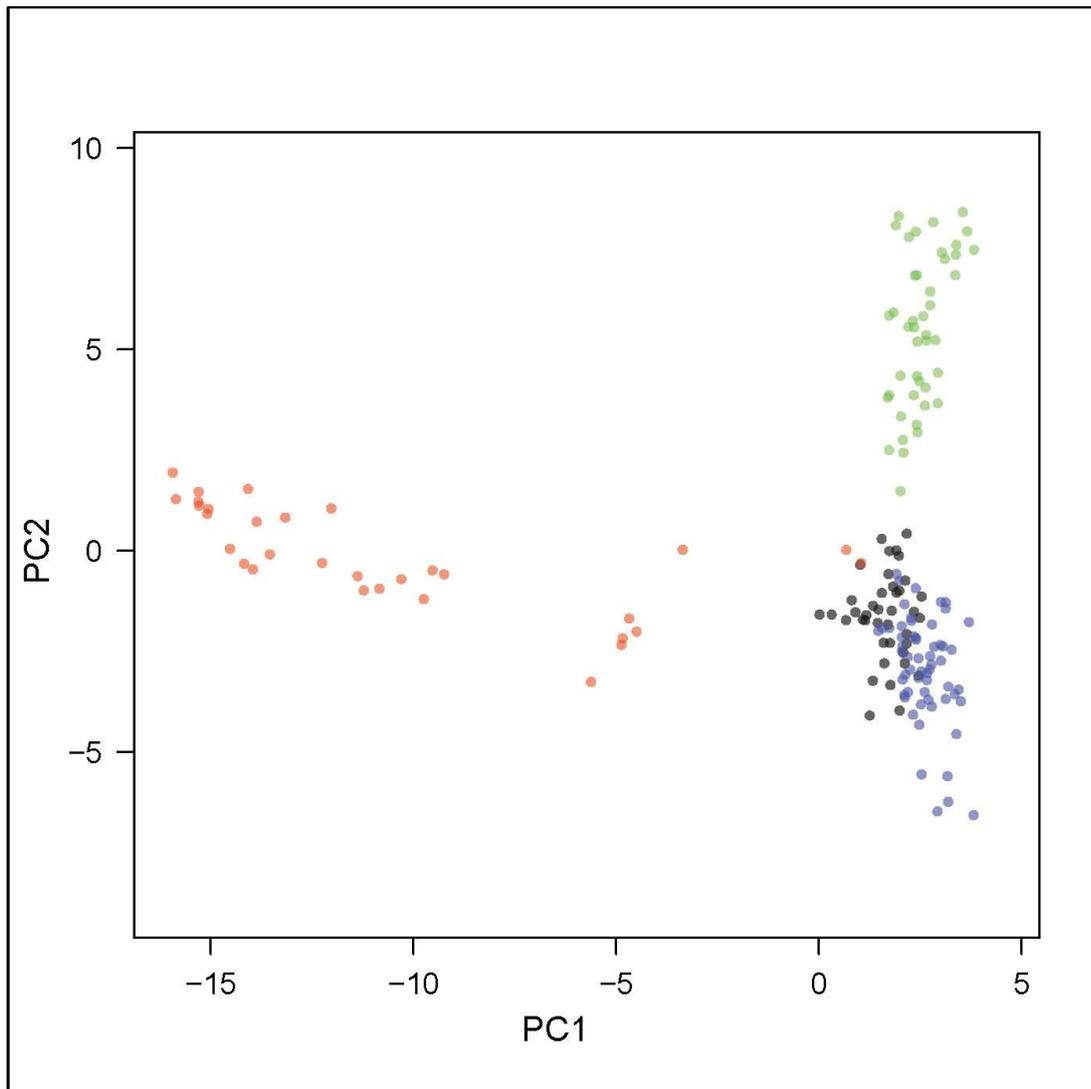


655 **Figure 3**



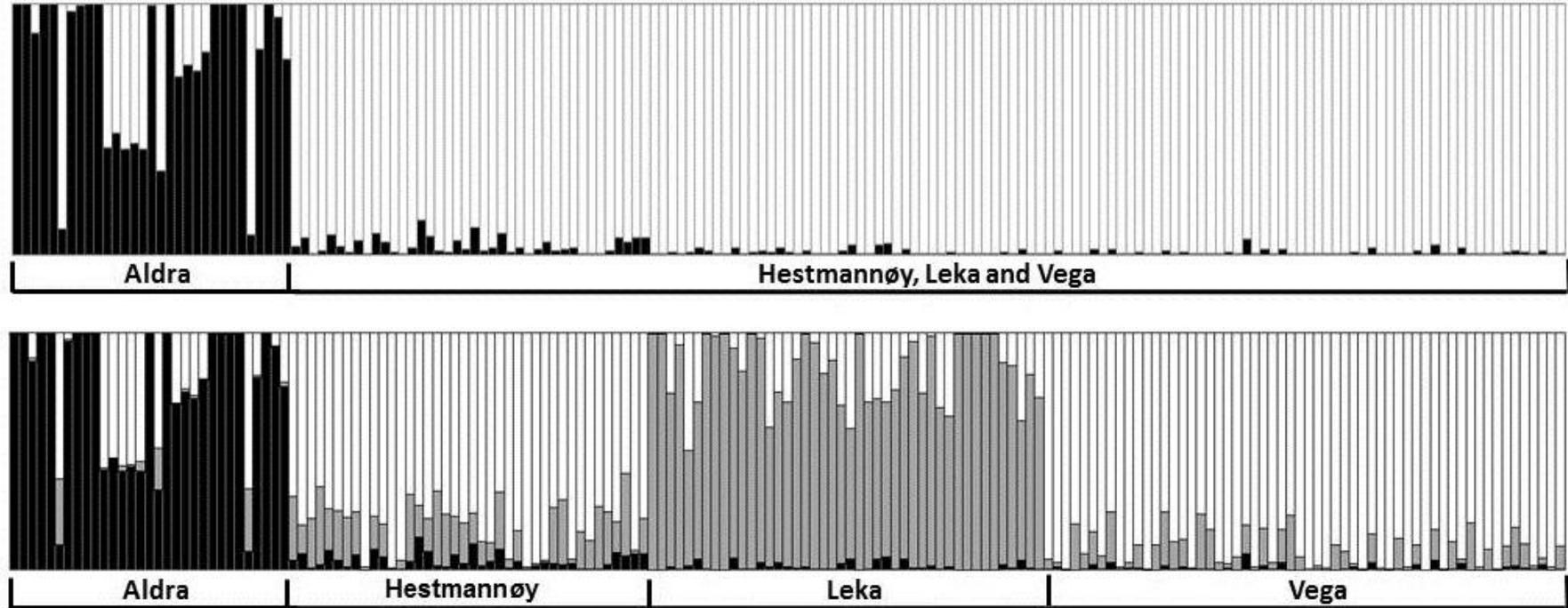
656

657

658 **Figure 4**

659

660 **Figure 5**



661

