1 **Supporting information on data analysis**

2 **Sampling, DNA extraction, quality assessment**

3 Sampling was conducted during the 2010/2011 summer season in the king penguin colony of

4 'La Baie du Marin' on Possession Island (46°25'S, 51°45'E) in the Crozet Archipelago. Blood (~

5 100 μL) was collected from the brachial vein of chicks hatched in the long-term monitored

6 area 'ANTAVIA', transferred to a filter paper (Whatman 113 ®), dried, and later frozen at -

7 20°C. Individuals were randomly selected along a 120m-axis at the periphery of the colony,

8 in order to maximise distance separation. A total of 140 individuals were chosen for

9 mitochondrial DNA Control Region analysis, and 8 of these were randomly selected for

10 restriction site-associated DNA (RAD) sequencing analysis. Total DNA were extracted from

11 the filter papers using a Phenol-Chloroform protocol or the Qiagen DNAase blood & tissue

12 kit according to manufacturer's instructions. After extraction, DNA quantity and quality were

13 tested in each sample by fluorimetric-based measurement (Qubit, Invitrogen) and gel

14 electrophoresis.

15

16 **MtDNA marker analysis**

17 Partial sequences of the Control Region (354 bp) were amplified and sequenced in 140

18 samples according to the protocol published in Heupkin *et al.* [1]. PCR products were Sanger-

19 sequenced in the ABI-LAB at the University of Oslo. Sequences were then manually edited

20 and aligned in Bioedit [2]. All new haplotype sequences have been uploaded to GenBank

21 (Accession number: KF530582-KF530720). Summary molecular statistics, demographic

22 parameters and the mismatch distribution of pairwise differences were calculated in DNAsp

23 v5 [3]. This dataset was used to infer the king penguin past demography employing the

24 Bayesian Skyride plot [4], where inferred population history is bounded by credibility

25 intervals that combine phylogenetic and coalescent uncertainties, as implemented in the

26 BEAST 1.7.4 package [5]. Analyses were performed on the Bioportal facility (now LifePortal)

27 running on the ABEL cluster, University of Oslo. A GTR+G+I substitution model was set for

28 the mitochondrial sequence. A relaxed uncorrelated log-normal clock prior was set for the

29 substitution rate to take into account fluctuations of the molecular clock along different

30   branches of the phylogeny; a log-normal priors with mean in the real space of 0.55

31   substitution/site/Myr respectively were set [6]. The Bayesian Skyline plot was set as

32   coalescent tree prior model. Convergence among three runs, with a MCMC length of 30

33   million generations for each parameter setting was checked. Effective sample size was

34   checked in Tracer 1.5 [7] and plots of population size change through time were drawn.

35

## RAD sequencing and genome-wide demographic inference

37   Eight king penguin individuals were pooled and genotyped by RAD sequencing [8] in one

38   library sequenced on an ILLUMINA HiSeq2000, yielding ca. 65 million 100-bp reads. All raw

39   sequence reads are available on GenBank at the Sequence Read Archive (Run Num.:

40   SRR942341).

41        After quality assessment, samples showing high molecular weight and highly

42   concentrated DNA were employed in next-generation sequencing (NGS) of RADtags [8]. The

43   following RADseq protocol was adopted: (i) approximately 100ng of genomic DNA per

44   sample were digested with the restriction enzyme *Sbf*I (NEB); (ii) each sample was then

45   ligated to a unique barcoded P1 adapter prior to pooling in a single library. The library was

46   then sheared by sonication, and gel electrophoresis of small library aliquots were run after

47   the first 5 cycles (30'' ON – 30'' OFF) and then every 2-3 cycles of sonication; (iii) the target

48   size range fraction (300-500 bp) was achieved after 8 cycles of sonication and was then

49   selected by gel electrophoresis and manual excision; (iv) before size selection on the gel,

50   sonicated libraries were concentrated to 25 µl by DNA capture on magnetic beads (beads

51   solution:DNA = 0.8:1), thus further reducing the carry-over of non-ligated P1 adapters; (v)

52   capture on magnetic beads using the same beads:DNA ratio (0.8:1) was then employed in all

53   following purification steps (after blunt-end repairing, poly-A tailing, P2 adapter ligation and

54   library enrichment by PCR); (vi) PCR amplification was performed in 8 x 12.5 µl aliquots

55   pooled after the amplification in order to reduce amplification bias on few loci due to

56   random drift; (vii) the library was then quantified by a fluorimetric-based method (Qubit,

57   Invitrogen) and molarity was checked on an Agilent Bioanalyzer chip (Invitrogen). A final

58   volume of 20 µl with a DNA concentration of 45 ng/µl was submitted for sequencing on an

59   ILLUMINA HiSeq2000 sequencer at the Norwegian Sequencing Centre, University of Oslo.

60    Raw reads were then processed using the scripts included in the Stacks package [9]

61    running on our server facility on the ABEL cluster, University of Oslo. Raw reads were quality

62    filtered and grouped according to individual barcodes. Then individual loci were retrieved

63    and SNPs were called by a maximum-likelihood function that excluded likely sequencing

64    errors. Several runs with different settings of read trimming parameter, quality thresholds,

65    mismatches allowed when building the individual and the population catalogs, were

66    performed to check for consistency of the results. The parameters setting used to build the

67    final catalog included: -t 95 and the default values for the quality checking when using

68    "process_radtags.pl"; -m 10, -n 7, -M 3 when running "denovo_map.pl". 101,115 loci with

69    50X average coverage were aligned in an unreferenced catalog. A table including all loci

70    matching the eight sequenced individuals was built using "export_sql.pl" Stacks script. This

71    table was further filtered by python scripts (available upon request) excluding loci with

72    missing data, with more than 2 alleles per individual, and deleveraged by Stacks algorithm.

73    Loci were then grouped according to the number of SNPs allowing from 0 to a

74    maximum of 6 substitutions per locus (0 to 6-SNP classes). Loci with 4-6 SNPs were then

75    directly checked through the catalog web-based interface provided by Stacks. Loci with more

76    than 2 SNPs in the last 5 base pairs or with observed heterozygosity higher than 0.6 were

77    blacklisted and removed from the table as likely sequencing errors or paralogous loci. Only

78    those loci hosting 1 single bi-allelic SNP were employed in AFs analysis in order to minimize

79    linkage among the data. Not having a reference genome, we could not exclude loci produced

80    by adjacent genomic regions or by the two sides of each restriction site. Custom python

81    scripts (available upon request) were employed to edit this 1-SNP dataset as a suitable input

82    file for downstream statistical analysis encoding SNPs as 0-2 when homozygote for the two

83    alleles respectively or 1 when heterozygote. On the other hand, loci in 2 to 6-SNP classes

84    were treated as short sequences and locus-by-locus edited using python script as NEXUS

85    format files each containing 16 sequences 95 bp long (two sequences per individual).

86    Minor allele frequency spectrum was calculated by functions available in the R

87    package "adegenet" [10] using loci included in the 1-SNP class. This information was then

88    passed to the python-based software ∂a∂i [11] that, using a diffusion approximation to the

89    allele frequency spectrum, allows demographic inference from genetic data testing

90    alternative demographic scenarios in a maximum-likelihood framework. A sudden growth in

population size was tested against the null hypothesis of constant population size using the "two_epoch" and the "snc" functions, respectively. Several runs of likelihood optimization were performed changing the extent of the search by the "fold" parameter in the "dadi.Misc.perturb_params" function. Optimized log-likelihood and Theta values were recorded. In order to calculate effective population size from Theta values produced by ∂a∂I, a total sequence length of 1,943,510 bp (95 bp X 20,458 loci used in this analysis) was used.

Functions included in the R package "ape" [12] and the R standard boxplot function [13] were used to calculate the joint mismatch distribution in the pairwise differences (from here onwards referred to as mismatch distribution density). Calculations were performed and plotted in each 2 to 6-SNP classes separately.

Different random combinations of 50-100 loci in 2 to 6-SNP classes were compared when inferring the past demography of the king penguin population using the coalescent-theory based multi-locus analysis implemented in BEAST 1.7.4. Linkage disequilibrium was tested in all subsets using Genepop [14] with the default setting in the web tool and the Bonferroni correction for multiple tests. The robustness of the approach was tested with respect to *i*) the number of SNPs per locus, *ii*) the different random selection of loci and *iii*) the number of loci included in the random selection: 50 loci in 2-SNP class (5 runs), 50 loci in 3-SNP class (5 runs), 50 loci in 4-6-SNP class (10 runs), 10 loci in 4-6-SNP class (1 run), 25 loci in 4-6-SNP class (1 run) and 100 loci in 4-6-SNP class (1 run). Three runs showing hints of multiple optima for the demographic function were discarded. Different settings of the parameters and priors have been explored in preliminary analyses, but the following was the definitive setting: (i) markers were unlinked concerning site substitution model, clock model and tree prior model; (ii) site substitution model was set as a HKY with empirical base frequency; (iii) a strict molecular clock was estimated for each marker with a uniform prior distribution bounded within 0.5 and 0.005 sub/s/Myr; (iv) the Extended Bayesian Skyline Plot (EBSP; [15]) was selected as tree prior model and ploidy of the markers was set accordingly. Fine tuning of operators did not improve our results as running the analyses with longer MCMC simulations; 200 million iterations were set as run length.

Mitochondrial Control Region data were included in the analysis in order to test the consistency of the information provided by the two genetic dataset (genomic and

121 mitochondrial) and to calibrate the genomic substitution rate using the mtDNA Control
122 Region substitution rate as estimated in the Adélie penguin. Site substitution and clock
123 models were set as in the analyses of the mitochondrial marker alone (see below). All
124 analyses were run on the Bioportal facility (now LifePortal) of the Abel cluster, University of
125 Oslo. Results were checked on Tracer 1.5 and plot of the EBSP data were drawn in R [13]. An
126 extensive study on Adélie penguin ancient DNA suggested a fast estimate (0.88 sub/s/Myr;
127 [16]) for the substitution rate of the mitochondrial Control Region. Further analyses
128 confirmed this high figure but it was downscaled to 0.55 sub/s/Myr [6]. We used the more
129 conservative 0.55 sub/s/Myr for our calibrated demography. A generation time of 11.49
130 years (Le Bohec *in prep*) was used to convert the population size estimates on the EBSP
131 (given by default in effective population size * generation time). We then plotted the
132 population trend for the last 35,000 years together with the trend of temperature anomalies
133 as inferred by the analysis of the EPICA Dome C ice core [17]. Concerning the calibration of
134 the mean genome-wide substitution rate: first, the mean of the median values in each SNP-
135 ratio class included in the EBSP analysis was calculated (4 to 6-SNP classes); then, a linear
136 regression was used to infer the substitution rate of those SNP classes excluded from the
137 EBSP analysis (0 to 3-SNP classes); finally, we calculated the mean genomic substitution rate
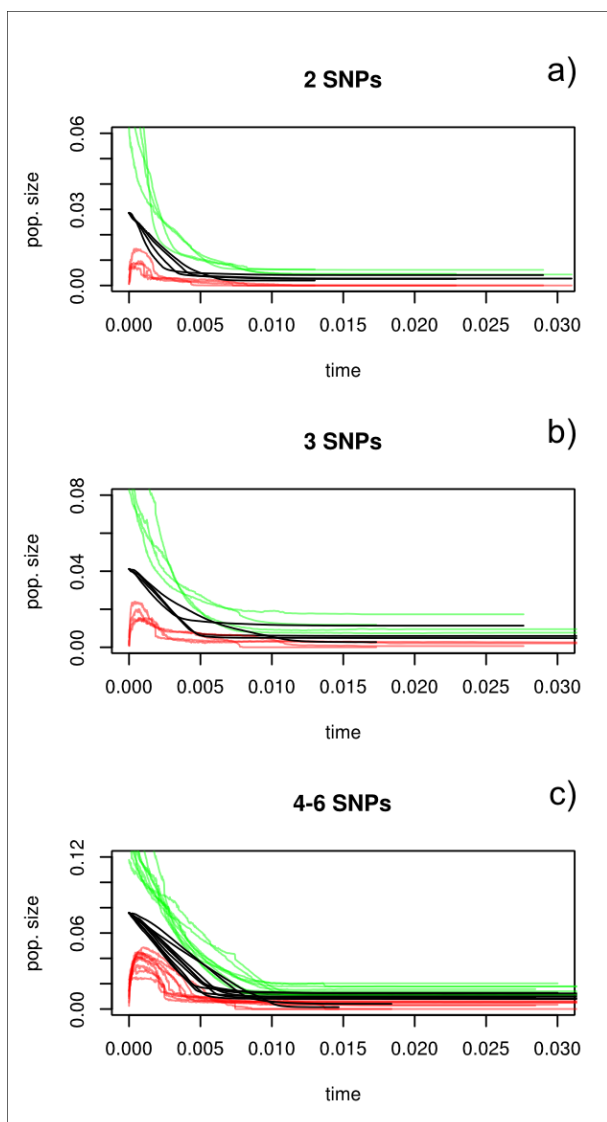138 weighting each SNP class accordingly with the frequency (number of loci) of each class.
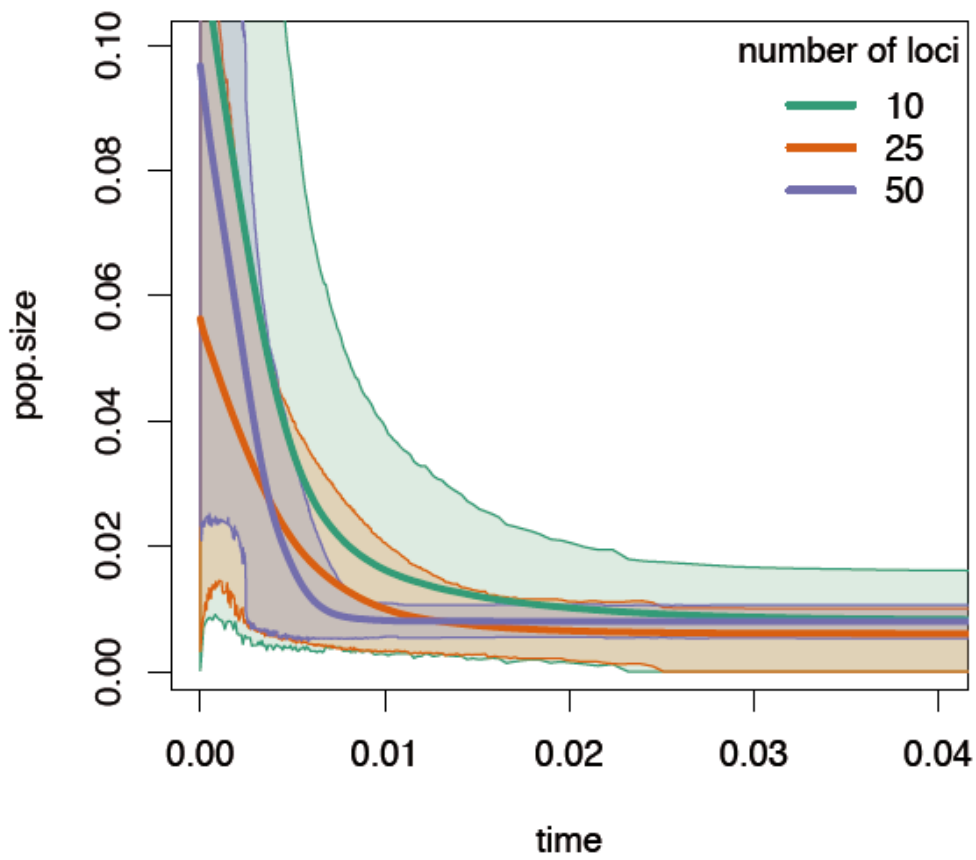
139

140 **References**

141 1. Heupink T. H., van den Hoff J., Lambert D. M. 2012 King penguin population on Macquarie Island
142    recovers ancient DNA diversity after heavy exploitation in historic times. *Biol. Lett.* **8**, 586–589.

143 2. Hall T. A.  1999 BioEdit: a user-friendly biological sequence alignment editor and analysis program for
144    Windows 95/98/NT. *Nucl. Acids Symp. Ser.* **41**, 95-98.

145 3. Librado P., Rozas J. 2009 DnaSP v5: a software for comprehensive analysis of DNA polymorphism data.
146    *Bioinformatics* **25**, 1451–1452.

147 4. Minin V. N., Bloomquist E. W. & Suchard M. A. 2008 Smooth Skyride through a Rough Skyline: Bayesian
148    Coalescent-Based Inference of Population Dynamics. *Mol Biol Evol* **25**, 1459–1471.

149 5. Drummond A. J. & Rambaut A 2007 BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evol*

150     *Biol* **7**, 214.

151   6.   Millar, C. D., Dodd, A., Anderson, J., Gibb, G. C., Ritchie, P. A., Baroni, C., Woodhams M. D., Hendy M. D.,
152        Lambert D. M. 2008 Mutation and evolutionary rates in Adélie penguins from the Antarctic. *PLoS*
153        *Genetics*, **4**, e1000209.

154   7.   Rambaut A. &  Drummond A. 2005 "Tracer version 1.3: a program for analyzing results from Bayesian
155        MCMC programs such as BEAST and MrBayes." Distributed by the authors at: http://evolove. zoo. ox. ac.
156        uk/software. html

157   8.   Baird, N. A., Etter, P. D., Atwood, T. S., Currey, M. C., Shiver, A. L., Lewis, Z. A., Selker, E. U., Cresko, W. A.
158        & Johnson, E. A. 2008 Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers. *PLoS*
159        *ONE* **3**, e3376. (doi:10.1371/journal.pone.0003376)

160   9.   Catchen J. M., Amores A., Hohenlohe P., Cresko W. & Postlethwait J. H. 2011 Stacks: Building and
161        Genotyping Loci De Novo From Short-Read Sequences. *G3* **1**, 171–182.

162   10.  Jombart T. & Ahmed I. 2011 adegenet 1.3-1: new tools for the analysis of genome-wide SNP data.
163        *Bioinformatics* **27**, 3070–3071.

164   11.  Gutenkunst R. N., Hernandez R. D., Williamson S. H. & Bustamante C. D. 2009 Inferring the Joint
165        Demographic History of Multiple Populations from Multidimensional SNP Frequency Data. *PLoS Genet* **5**,
166        e1000695.

167   12.  Paradis E., Claude J. & Strimmer K. 2004 APE: Analyses of Phylogenetics and Evolution in R language.
168        *Bioinformatics* **20**, 289–290.

169   13.  R Development Core Team 2011 R Foundation for Statistical Computing, Vienna, Austria URL
170        http://www.R-project.org.

171   14.  Raymond M. & Rousset F. 1995 GENEPOP (version 1.2): population genetics software for exact tests and
172        ecumenicism. J. Heredity, 86, 248-249

173   15.  Heled J. & Drummond A.J. 2008 Bayesian inference of population size history from multiple loci. *BMC*
174        *Evol Biol* **8**, 289.

175   16.  Lambert, D. M., Ritchie, P. A., Millar, C. D., Holland, B., Drummond, A. J., & Baroni, C. 2002 Rates of
176        evolution in ancient DNA from Adélie penguins. *Science*, **295**, 2270-2273.

177   17.  Jouzel, J., Masson-Delmotte, V., Cattani, O., Dreyfus, G., Falourd, S., Hoffmann, G., Minster B., Nouet J.,
178        Barnola J. M., Chappellaz J. *et al* 2007 Orbital and millennial Antarctic climate variability over the past
179        800,000 years. *Science*, **317**, 793-796.

180    **Figure S1**. Demographic reconstructions of the Crozet king penguin colony employing the

181    Extended Bayesian Skyline Plot analysis. Consistency in the pattern inferred is compared

182    among different data selections including 50 loci chosen at random from different classes of

183    variation: a) 2 SNPs, 4 independent datasets; b) 3 SNPs, 4 independent datasets; c) 4-6 SNPs,

184    9 independent datasets. In order to facilitate comparison of uncalibrated EBSP

185    reconstruction and solely for visualization purpose, all runs were scaled to have the same

186    value of the demographic function at t0. To do so, we divided all demographic estimates by

187    the ratio $Ki = Ni(t0)/max(Ni(t0))$, where $Ni(t0)$ is the median value from the posterior

188    distribution of the demographic function at t0 for each run. Correspondingly, the time

189    intervals of each run were multiplied by the same ratio Ki, reflecting the assumption that the

190    actual population size is the same across all runs. Median (black), 95% HPD lower (red) and

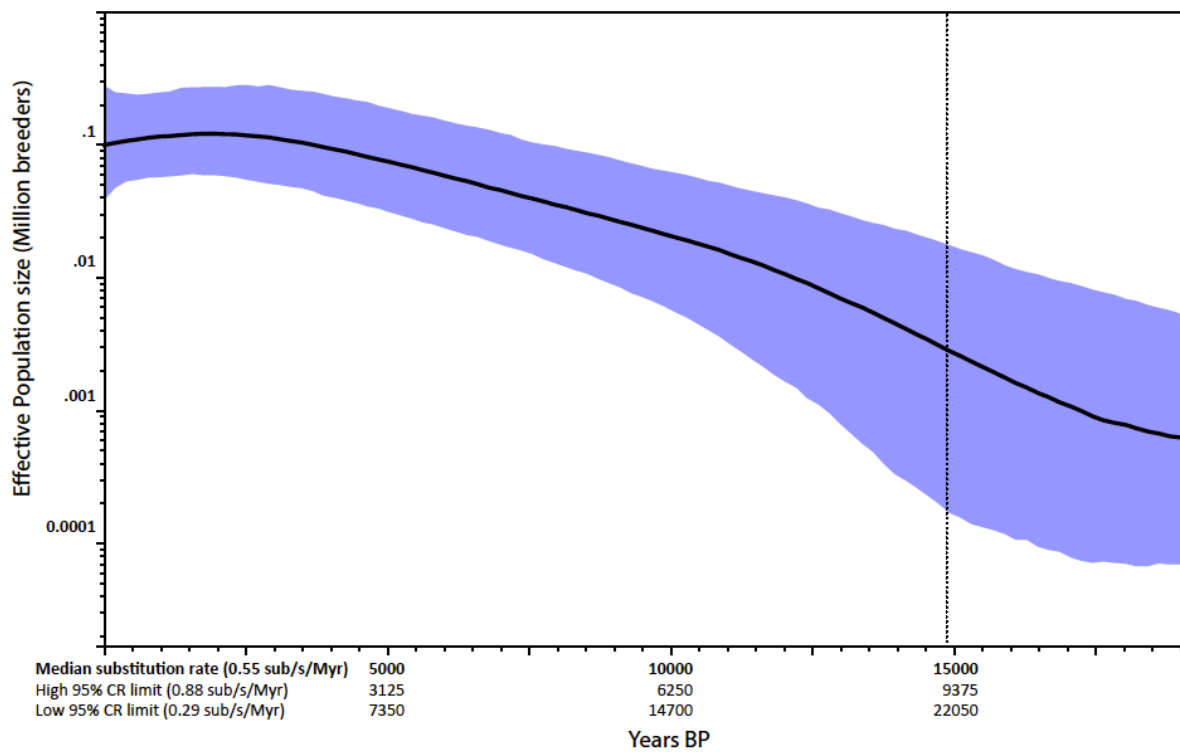191    upper (green) values are reported.



192

193 **Figure S2.** Demographic reconstructions of the Crozet king penguin colony employing the
194 Extended Bayesian Skyline Plot analysis. Consistency of the inference is compared across
195 three nested datasets (10, 25 and 50 loci) randomly selected from the 4-6-SNP class.
196 Population size and time are unscaled. Median (solid line) and 95% upper and lower (filled
197 areas) values are reported.



198

199

200

201

202

203

204

**Figure S3.** Bayesian Skyride Plot inferred from mitochondrial Control Region data describing
population trend through time. Time is scaled according to median (solid black line), 95%
upper and lower credibility region (filled blue area), as estimated in Millar *et al* [6].