# Programmed DNA elimination of germline development genes in songbirds

Cormac M. Kinsella[1,2]*, Francisco J. Ruiz-Ruano[3]*, Anne-Marie Dion-Côté[1,4], Alexander J. Charles[5], Toni I. Gossmann[5], Josefa Cabrero[3], Dennis Kappei[6], Nicola Hemmings[5], Mirre J. P. Simons[5], Juan P. M. Camacho[3], Wolfgang Forstmeier[7], Alexander Suh[1]

[1]Department of Evolutionary Biology, Evolutionary Biology Centre (EBC), Science for Life Laboratory, Uppsala University, SE-752 36, Uppsala, Sweden.

[2]Laboratory of Experimental Virology, Department of Medical Microbiology, Amsterdam UMC, University of Amsterdam, 1105 AZ, Amsterdam, The Netherlands.

[3]Department of Genetics, University of Granada, E-18071, Granada, Spain.

[4]Department of Molecular Biology & Genetics, Cornell University, NY 14853, Ithaca, United States.

[5]Department of Animal and Plant Sciences, University of Sheffield, S10 2TN, Sheffield, United Kingdom.

[6]Cancer Science Institute of Singapore, National University of Singapore, 117599, Singapore.

[7]Max Planck Institute for Ornithology, D-82319, Seewiesen, Germany.

*These authors contributed equally to this work (alphabetical order).

Correspondence and requests for materials should be addressed to F.J.R.R. (email: fjruizruano@ugr.es) and A.S. (email: alexander.suh@ebc.uu.se).

24 **Summary**

25 Genomes can vary within individual organisms. Programmed DNA elimination leads to dramatic

26 changes in genome organisation during the germline–soma differentiation of ciliates[1], lampreys[2],

27 nematodes[3,4], and various other eukaryotes[5]. A particularly remarkable example of tissue-specific

28 genome differentiation is the germline-restricted chromosome (GRC) in the zebra finch which is

29 consistently absent from somatic cells[6]. Although the zebra finch is an important animal model

30 system[7], molecular evidence from its large GRC (>150 megabases) is limited to a short intergenic

31 region[8] and a single mRNA[9]. Here, we combined cytogenetic, genomic, transcriptomic, and

32 proteomic evidence to resolve the evolutionary origin and functional significance of the GRC.

33 First, by generating tissue-specific *de-novo* linked-read genome assemblies and re-sequencing

34 two additional germline and soma samples, we found that the GRC contains at least 115 genes

35 which are paralogous to single-copy genes on 18 autosomes and the Z chromosome. We detected

36 an amplification of ≥38 GRC-linked genes into high copy numbers (up to 185 copies) but,

37 surprisingly, no enrichment of transposable elements on the GRC. Second, transcriptome and

38 proteome data provided evidence for functional expression of GRC genes at the RNA and protein

39 levels in testes and ovaries. Interestingly, the GRC is enriched for genes with highly expressed

40 orthologs in chicken gonads and gene ontologies involved in female gonad development. Third,

41 we detected evolutionary strata of GRC-linked genes. Genes such as *bicc1* and *trim71* have

42 resided on the GRC for tens of millions of years, whereas dozens have become GRC-linked very

43 recently. The GRC is thus likely widespread in songbirds (half of all bird species) and its rapid

44 evolution may have contributed to their diversification. Together, our results demonstrate a

45 highly dynamic evolutionary history of the songbird GRC leading to dramatic germline–soma

46 genome differences as a novel mechanism to minimize genetic conflict between germline and

47 soma.

48  **Text**

49  Not all cells of an organism must contain the same genome. Some eukaryotes exhibit dramatic

50  differences between their germline and somatic genomes, resulting from programmed DNA

51  elimination of chromosomes or fragments thereof during germline–soma differentiation[5]. Here

52  we present the first comprehensive analyses of a germline-restricted chromosome (GRC). The

53  zebra finch (*Taeniopygia guttata*) GRC is the largest chromosome of this songbird[6] and likely

54  comprises >10% of the genome (>150 megabases)[7,10]. Cytogenetic evidence suggests the GRC is

55  inherited through the female germline, expelled late during spermatogenesis, and eliminated from

56  the soma during early embryo development[6,11]. Previous analyses of a 19-kb intergenic region

57  suggested that the GRC contains sequences with high similarity to regular chromosomes ('A

58  chromosomes')[8].

59

60  In order to reliably identify sequences as GRC-linked, we used a single-molecule sequencing

61  technology not applied in birds before that permits reconstruction of long haplotypes through

62  linked reads[12]. We generated separate haplotype-resolved *de-novo* genome assemblies for the

63  germline and soma of a male zebra finch (testis and liver; 'Seewiesen'; Supplementary Table 1).

64  We further used the linked-read data to compare read coverage and haplotype barcode data in

65  relation to the zebra finch somatic reference genome ('taeGut2')[7], allowing us to identify

66  sequences that are shared, amplified, or unique to the germline genome in a fashion similar to

67  recent studies on cancer aneuploidies[13]. We also re-sequenced the germline and soma from two

68  male zebra finches from another population ('Spain'; testis and muscle) using short reads.

69

3

70    We first established the presence of the GRC in the three germline samples. Cytogenetic analysis

71    using fluorescence *in-situ* hybridization (FISH) with a new GRC probe showed that the GRC is

72    present exclusively in the germline and eliminated during spermatogenesis as hypothesised (Fig.

73    1a-b, Extended Data Fig. 1)[6,11]. We compared germline/soma sequencing coverage by mapping

74    reads from all three sampled zebra finches onto the reference genome assembly (regular 'A

75    chromosomes'), revealing consistently germline-increased coverage for single-copy regions,

76    similar to programmed DNA elimination of short genome fragments in lampreys[2] (Fig. 1c-d). A

77    total of 92 regions (41 with >10 kb length) on 13 chromosomes exhibit >4-fold increased

78    germline coverage in Seewiesen relative to the soma (Fig. 1e, Supplementary Table 2). Such a

79    conservative coverage cut-off provides high confidence in true GRC-amplified regions. We

80    obtained nearly identical confirmatory results with another sequencing technology using the

81    Spanish birds (Fig. 1f). Notably, the largest block of testis-increased coverage spans nearly 1 Mb

82    on chromosome 1 and overlaps with the previously[8] FISH-verified intergenic region 27L4 (Fig.

83    1e-f).

84

85    Our linked-read and re-sequencing approach allowed us to determine the sequence content of the

86    GRC. The GRC is effectively a non-recombining chromosome as it recombines with itself after

87    duplication, probably to ensure stable inheritance during female meiosis[8]. We predicted that the

88    GRC would be highly enriched in repetitive elements, similar to the female-specific avian W

89    chromosome (repeat density >50%)[14]. Surprisingly, neither assembly-based nor read-based repeat

90    quantifications detected a significant enrichment in transposable elements or satellite repeats in

91    the germline samples relative to the soma samples (Supplementary Text, Supplementary Table 3).

92    Instead, most germline coverage peaks lie in single-copy regions of the reference genome

4

93  including 38 genes (Fig. 1e-f, Table 1), suggesting that these peaks stem from highly similar

94  GRC-amplified paralogs with high copy numbers (up to 185 copies per gene; Supplementary

95  Table 4). GRC linkage of these regions is further supported by sharing of linked-read barcodes

96  between different amplified chromosomal regions in germline but not soma (Fig. 1g-h),

97  suggesting that these regions reside on the same haplotype (Extended Data Fig. 2). We

98  additionally identified 245 GRC-linked genes through germline-specific single-nucleotide

99  polymorphisms (SNPs) present in read mapping of all three germline samples onto zebra finch

100 reference genes (up to 402 SNPs per gene; Supplementary Table 4). As a control, we used the

101 same methodology to screen for soma-specific SNPs and found no such genes. We

102 conservatively consider the 38 GRC-amplified genes and those with at least 5 germline-specific

103 SNPs as our highest-confidence set (Table 1). We also identified GRC-linked genes using

104 germline–soma assembly subtraction and coverage differences (Fig. 1i); however, all were

105 already found via coverage or SNP evidence (Table 1). Together with the *napa* gene recently

106 identified in transcriptomes (Fig. 1j)[9], our complementary approaches yielded 115 high-

107 confidence GRC-linked genes with paralogs located on 18 autosomes and the Z chromosome

108 (Table 1; all 267 GRC genes in Supplementary Table 4).

109

110 We next tested whether the GRC is functional and thus probably physiologically important using

111 transcriptomics and proteomics. We sequenced RNA from the same tissues of the two Spanish

112 birds used for genome re-sequencing and combined these with published testis and ovary RNA-

113 seq data from North American domesticated zebra finches[9,15]. Among the 115 high-confidence

114 genes, 6 and 32 were transcribed in testes and ovaries, respectively (Table 1). Note, these are

115 only genes where we could reliably separate reads from GRC-linked and A-chromosomal

116    paralogs by GRC-specific SNPs in the transcripts (Fig. 2a-b, Extended Data Fig. 3,

117    Supplementary Table 5). We next verified translation of GRC-linked genes through protein mass

118    spectrometry data for 7 testes and 2 ovaries from another population ('Sheffield'). From 83 genes

119    with GRC-specific amino acid changes, we identified peptides from 5 GRC-linked genes in testes

120    and ovaries (Fig. 2c-d, Extended Data Fig. 4, Table 1, Supplementary Table 6). We hence

121    established that many GRC-linked genes are transcribed and translated in adult male and female

122    gonads, extending previous RNA evidence for a single gene[9] and refuting the hypothesis from

123    cytogenetic studies that the GRC is silenced in the male germline[16,17]. Instead, we hypothesise

124    that the GRC has important functions during germline development, which is supported by a

125    significant enrichment in gene ontology terms related to reproductive developmental processes

126    among GRC-linked genes (Fig. 2e, Supplementary Table 7). We further found that the GRC is

127    significantly enriched in genes that are also germline-expressed in GRC-lacking species with

128    available RNA expression data from many tissues[18] (Fig. 2f, Supplementary Table 8).

129    Specifically, we found that 22 and 6 out of 65 chicken orthologs of high-confidence GRC-linked

130    genes are most strongly expressed in chicken testis and ovary, respectively.

131

132    The observation that all identified GRC-linked genes have A-chromosomal paralogs allowed us

133    to decipher the evolutionary origins of the GRC. We utilised phylogenies of GRC-linked genes

134    and their A-chromosomal paralogs to infer when these genes copied to the GRC, similarly to the

135    inference of evolutionary strata of sex chromosome differentiation[19]. First, the phylogeny of the

136    intergenic 27L4 locus of our germline samples and a previous GRC sequence[8] demonstrated

137    stable inheritance among the sampled zebra finch populations (Fig. 3a). Second, 37 gene trees of

138    GRC-linked genes with germline-specific SNPs and available somatic genome data from other

139     birds identify at least five evolutionary strata (Fig. 3b-f, Extended Data Fig. 5, Table 1), with all

140     but stratum 3 containing expressed genes (*cf.* Fig. 2a-d). Stratum 1 emerged during early songbird

141     diversification, stratum 2 before the diversification of estrildid finches, and stratum 3 within

142     estrildid finches (Fig. 3g). The presence of at least 7 genes in these three strata implies that the

143     GRC is tens of millions of years old and likely present across songbirds (Extended Data Fig. 5),

144     in line with a recent cytogenetics preprint[20]. Notably, stratum 4 is specific to the zebra finch

145     species and stratum 5 to the Australian zebra finch subspecies (Fig. 3g), suggesting piecemeal

146     addition of genes from 18 autosomes and the Z chromosome over millions of years of GRC

147     evolution (Fig. 3h). The long-term residence of expressed genes on the GRC implies that they

148     have been under selection, such as *bicc1* and *trim71* on GRC stratum 1 whose human orthologs

149     are important for embryonic cell differentiation[21]. Additionally, we detected evidence for

150     purifying selection on GRC-linked genes from older and younger strata using ratios of non-

151     synonymous to synonymous substitutions (dN/dS; Supplementary Table 9), again implying that

152     the GRC is an important chromosome with a long evolutionary history.

153

154     Here we provided the first evidence for the origin and functional significance of a GRC. Notably,

155     our analyses suggest that the GRC emerged during early songbird evolution and we predict it to

156     be present in half of all bird species. The species-specific addition of dozens of genes on stratum

157     5 implies that the rapidly evolving GRC likely contributed to reproductive isolation during the

158     massive diversification of songbirds[22]. It was previously hypothesised that GRCs are formerly

159     parasitic B chromosomes that became stably inherited after acquiring essential functions for the

160     host[23,24]. Our evidence for an enrichment of germline-expressed genes on the zebra finch GRC is

161     reminiscent of nematodes and lampreys where short genome fragments containing similar genes

162 are eliminated during germline–soma differentiation[2-4]. All these cases constitute extreme

163 mechanisms of gene regulation through germline–soma gene removal rather than transcriptional

164 repression[3,5,10]. Consequently, we hypothesise that the GRC became indispensable for its host by

165 the acquisition of germline development genes. The aggregation of developmental genes on a

166 single eliminated chromosome constitutes a novel mechanism to ensure germline-specific gene

167 expression in multicellular organisms. This may allow adaptation to germline-specific functions

168 free of deleterious effects on the soma which would otherwise arise from antagonistic pleiotropy.

169

170 **References (max. 30 references)**

171 1    Chen, X. *et al.* The architecture of a scrambled genome reveals massive levels of genomic
172      rearrangement during development. *Cell* **158**, 1187-1198 (2014).
173 2    Smith, J. J. *et al.* The sea lamprey germline genome provides insights into programmed
174      genome rearrangement and vertebrate evolution. *Nat. Genet.* **50**, 270-277 (2018).
175 3    Wang, J. *et al.* Silencing of germline-expressed genes by DNA elimination in somatic
176      cells. *Dev. Cell* **23**, 1072-1080 (2012).
177 4    Wang, J. *et al.* Comparative genome analysis of programmed DNA elimination in
178      nematodes. *Genome Res.* **27**, 2001-2014 (2017).
179 5    Wang, J. & Davis, R. E. Programmed DNA elimination in multicellular organisms. *Curr.*
180      *Opin. Genet. Dev.* **27**, 26-34 (2014).
181 6    Pigozzi, M. I. & Solari, A. J. Germ cell restriction and regular transmission of an
182      accessory chromosome that mimics a sex body in the zebra finch, *Taeniopygia guttata.*
183      *Chromosome Res.* **6**, 105-113 (1998).
184 7    Warren, W. C. *et al.* The genome of a songbird. *Nature* **464**, 757–762 (2010).
185 8    Itoh, Y., Kampf, K., Pigozzi, M. I. & Arnold, A. P. Molecular cloning and
186      characterization of the germline-restricted chromosome sequence in the zebra finch.
187      *Chromosoma* **118**, 527-536 (2009).
188 9    Biederman, M. K. *et al.* Discovery of the first germline-restricted gene by subtractive
189      transcriptomic analysis in the zebra finch, *Taeniopygia guttata. Curr. Biol.* **28**, 1620-1627
190      (2018).
191 10   Smith, J. J. Programmed DNA elimination: keeping germline genes in their place. *Curr.*
192      *Biol.* **28**, R601-R603 (2018).
193 11   Pigozzi, M. I. & Solari, A. J. The germ-line-restricted chromosome in the zebra finch:
194      recombination in females and elimination in males. *Chromosoma* **114**, 403-409 (2005).
195 12   Weisenfeld, N. I., Kumar, V., Shah, P., Church, D. M. & Jaffe, D. B. Direct determination
196      of diploid genome sequences. *Genome Res.* **27**, 757-767 (2017).

197  13      Bell, J. M. *et al.* Chromosome-scale mega-haplotypes enable digital karyotyping of cancer
198          aneuploidy. *Nucleic Acids Res.* **45**, e162-e162 (2017).
199  14      Kapusta, A. & Suh, A. Evolution of bird genomes—a transposon's-eye view. *Ann. N. Y.*
200          *Acad. Sci.* **1389**, 164–185 (2017).
201  15      Singhal, S. *et al.* Stable recombination hotspots in birds. *Science* **350**, 928-932 (2015).
202  16      del Priore, L. & Pigozzi, M. I. Histone modifications related to chromosome silencing and
203          elimination during male meiosis in Bengalese finch. *Chromosoma* **123**, 293-302 (2014).
204  17      Goday, C. & Pigozzi, M. I. Heterochromatin and histone modifications in the germline-
205          restricted chromosome of the zebra finch undergoing elimination during spermatogenesis.
206          *Chromosoma* **119**, 325-336 (2010).
207  18      Marin, R. *et al.* Convergent origination of a *Drosophila*-like dosage compensation
208          mechanism in a reptile lineage. *Genome Res.* **27**, 1974-1987 (2017).
209  19      Lahn, B. T. & Page, D. C. Four evolutionary strata on the human X chromosome. *Science*
210          **286**, 964-967 (1999).
211  20      Torgasheva, A. A. *et al.* Germline-restricted chromosome (GRC) is widespread among
212          songbirds. *bioRxiv* **doi:10.1101/414276** (2018).
213  21      Uhlén, M. *et al.* Tissue-based map of the human proteome. *Science* **347** (2015).
214  22      Moyle, R. G. *et al.* Tectonic collision and uplift of Wallacea triggered the global songbird
215          radiation. *Nat. Commun.* **7**, 12709 (2016).
216  23      Camacho, J. P. M. B chromosomes. in *The Evolution of the Genome* (ed T. Ryan
217          Gregory) 223-286 (Elsevier Academic Press, 2005).
218  24      Camacho, J. P. M., Sharbel, T. F. & Beukeboom, L. W. B-chromosome evolution. *Philos.*
219          *Trans. R. Soc. B* **355**, 163-178 (2000).
220  25      Hooper, D. M. & Price, T. D. Rates of karyotypic evolution in Estrildid finches differ
221          between island and continental clades. *Evolution* **69**, 890-903 (2015).

222

245

246   **Author Contributions** Conceptualization: W.F., A.S., J.P.M.C., F.J.R.R., C.M.K., A.M.D.C.;

247   cytogenetics: J.P.M.C., F.J.R.R., J.C.; genomics: A.S., C.M.K., F.J.R.R., A.M.D.C.;

248   transcriptomics: F.J.R.R.; proteomics: T.I.G., A.J.C., D.K., M.S., N.H.; gene enrichment: C.M.K.,

249   W.F.; phylogenies: F.J.R.R., A.S., C.M.K.; manuscript writing: A.S. with input from all authors;

250   supervision: A.S., J.P.M.C. All authors read and approved the manuscript.

251

252   **Author Information** The authors declare no competing financial interests. Correspondence and

253   requests for materials should be addressed to F.J.R.R. (email: fjruizruano@ugr.es) and A.S.

254   (alexander.suh@ebc.uu.se).

255

256   **Tables and Figures**

257

258 **Table 1 | The 115 high-confidence genes on GRC with information on their A-chromosomal origin in the**
259 **reference genome taeGut2, number of testis-specific SNPs, methods supporting their GRC linkage,**
260 **testis/ovary RNA expression of the GRC paralog, testis/ovary protein expression of the GRC paralog, and**
261 **evolutionary stratum on the GRC.**

| Gene symbol | Chr. | Start | End | SNPs | Method | RNA evidence | Protein evidence | GRC stratum |
|---|---|---|---|---|---|---|---|---|
| AAGAB | 10 | 19608548 | 19634367 | 10 | SNPs | | | S5 |
| ADGRL2 | 8 | 14047115 | 14171612 | 10 | SNPs | | | |
| ADGRL3 | 4 | 14919933 | 15404594 | 8 | SNPs | ovary | | |
| AKIRIN2 | 3 | 78683482 | 78688947 | 6 | SNPs | ovary | | S5 |
| ALDH18A1 | 6 | 36280145 | 36301392 | 17 | SNPs | | | S4 |
| ALG13 | 4A | 18474239 | 18501426 | 19 | SNPs | ovary | | |
| ARMC6 | 28 | 4942046 | 4946063 | 5 | SNPs | | | |
| ATP2A2 | 15 | 2841010 | 2879975 | 8 | SNPs | | | |
| BICC1 | 6 | 6355408 | 6434911 | 402 | SNPs | ovary | | S1 |
| BMP15 | 4A | 15596686 | 15598225 | 29 | SNPs, coverage | ovary | | S5 |
| BMPR1B | 4 | 18997710 | 19024248 | 47 | SNPs, coverage | | | S5 |
| CCND3 | 26_random | | | 14 | SNPs | | | |
| CD164 | 3 | 69169111 | 69174605 | 38 | SNPs, coverage | ovary | | |
| COPS2 | 10 | 10200701 | 10222248 | 1 | SNPs, coverage | ovary | | |
| CPEB1 | 10 | 3114181 | 3137661 | 114 | SNPs | ovary | | |
| CSNK1A1L | Un | 135422201 | 135425792 | NA | coverage | | | |
| CXCL14 | 13 | 9423543 | 9433139 | 12 | SNPs | | | S5 |
| DDX49 | 28 | 4913058 | 4918451 | 5 | SNPs | ovary | | |
| DIS3L | 10 | 19097281 | 19112154 | 13 | SNPs | ovary | | S5 |
| DNAAF5 | 14 | 13758049 | 13780402 | NA | coverage | | | |
| DNAH5 | 2 | 81235805 | 81361091 | 7 | SNPs | | | |
| DPH6 | 5 | 31543945 | 31606965 | 13 | SNPs, coverage | | | |
| EFNB1 | 4A | 5764021 | 5807953 | 86 | SNPs | ovary | | S5 |
| ELAVL4 | 8 | 21034240 | 21098310 | 364 | SNPs | ovary | | |

11

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| EPPK1 | Un | | | 52 | SNPs | | | |
| FBXO16 | 3 | 112541865 | 112568948 | 6 | SNPs | | | |
| FEM1B | 10 | 19886491 | 19891616 | 9 | SNPs | ovary | | S5 |
| FIG4 | 3 | 69023384 | 69073678 | 17 | SNPs | | | S5 |
| FRS3 | 26_random | | | 42 | SNPs, coverage | | | S5 |
| GBE1 | 1 | 105820640 | 105934310 | 4 | SNPs, coverage | | | |
| INTS9 | 3 | 112259951 | 112313512 | NA | coverage | | | |
| LIAS | 4 | 48132714 | 48139736 | 42 | SNPs | | | S2 |
| LIN54 | 4 | 13615974 | 13637371 | 17 | SNPs | | | |
| LINC02027 | 1 | 106086596 | 106087033 | NA | coverage | | | |
| LMBRD2 | Z | 41646446 | 41665840 | NA | coverage | | | |
| LOC100223190 | Z | 69149414 | 69156994 | 41 | SNPs | | | |
| LOC100224235 | Un | | | 5 | SNPs | | | S5 |
| LOC100225322 | 1A | 47543094 | 47544622 | 6 | SNPs | ovary | | |
| LOC100227189 | Un | 150797142 | 150801997 | NA | coverage | | | |
| LOC100228170 | Un | 55540047 | 55541360 | NA | coverage | | | |
| LOC101233087 | Z | 47991391 | 47994344 | 7 | SNPs | | | |
| LOC101233688 | 5 | 937818 | 939059 | 5 | SNPs | | | S5 |
| LOC101233767 | 18 | 8034939 | 8038005 | 11 | SNPs | | | |
| LOC101233800 | Un | | | 16 | SNPs | | | |
| LOC101234253 | 10 | 19184028 | 19186114 | 7 | SNPs | ovary | | S5 |
| LOC105758464 | 23 | 46808 | 60360 | 14 | SNPs | | | S5 |
| LOC105758894 | 26_random | | | 5 | SNPs | | | |

12

| LOC105758976 | 2 | 34301994 | 34306899 | 16 | SNPs | | | |
| LOC105759101 | 3 | 76396180 | 76401262 | 21 | SNPs | | | |
| LOC105759167 | 4A | 15573874 | 15574621 | 5 | SNPs | | | |
| LOC105759195 | 4 | 14453003 | 14473747 | 18 | SNPs | | | |
| LOC105759199 | 4 | 20714525 | 20720872 | 11 | SNPs | | | |
| LOC105759260 | 5 | 1874731 | 1886007 | 32 | SNPs | | | S5 |
| LOC105759646 | Un | | | 7 | SNPs | | | |
| LOC105759655 | Un | | | 8 | SNPs | | | |
| LOC105759660 | Un | | | 18 | SNPs | | | |
| LOC105759665 | Un | | | 5 | SNPs | | | |
| LOC105759692 | Un | | | 12 | SNPs | | | |
| LOC105759919 | Un | | | 8 | SNPs | | | |
| LOC105760011 | Un | | | 7 | SNPs | | | |
| LOC105760123 | Un | | | 18 | SNPs | | | |
| LOC105760228 | Un | | | 14 | SNPs | | | |
| LOC105760286 | Un | | | 18 | SNPs | | | |
| LOC105760461 | Un | | | 10 | SNPs | | | |
| LOC105760874 | Z | 60949696 | 60953194 | 19 | SNPs | testis | | |
| LOC105760936 | 16_random | | | 12 | SNPs | | | |
| LUC7L3 | Un | 35019850 | 35021569 | NA | coverage | | | |
| MED20 | 26_random | 110500 | 113183 | 28 | SNPs, coverage | | | S5 |

13

| MSH4 | 8 | 27964612 | 27983306 | 30 | SNPs | | | S4 |
|---|---|---|---|---|---|---|---|---|
| NAPA | NA | | | NA | Biederman et al. 2018 | | both | |
| NEUROG1 | 13 | 9450787 | 9451086 | 6 | SNPs | | | |
| NFYA | 26 | 4725655 | 4735626 | 7 | SNPs | | | S5 |
| NRBP2 | 2 | 156379345 | 156398225 | 48 | SNPs | | | |
| PCSK4 | 28 | 4059367 | 4063775 | 21 | SNPs | | | |
| PGC | 26_random | | | 24 | SNPs | | | |
| PHKA1 | 4A | 15562688 | 15593666 | 16 | SNPs | | | |
| PIM1 | 26 | 603349 | 607242 | 50 | SNPs | testis | | |
| PIM3 | 1A | 18426716 | 18430551 | 81 | SNPs | ovary | | |
| PMM1 | 1A | 49038672 | 49047011 | NA | coverage | | | |
| PRDM1 | 3 | 70624594 | 70644625 | 12 | SNPs | | | |
| PRKAR1A | 18 | 2200317 | 2211579 | NA | coverage | | | |
| PRKAR1B | 14 | 13784578 | 13872733 | NA | coverage | | | |
| PRPSAP1 | 18 | 8008870 | 8033058 | 7 | SNPs, coverage | ovary | | S5 |
| PSIP1 | Z | 59887174 | 59919902 | 57 | SNPs, coverage | ovary | | S3 |
| PUF60 | 2 | 156354670 | 156376091 | 63 | SNPs | ovary | | |
| RFC1 | 4 | 48169638 | 48202709 | 77 | SNPs | ovary | | S2 |
| RNF157 | 18 | 8048721 | 8062403 | NA | coverage | | | |
| RNF17 | 1 | 45827734 | 45870640 | 69 | SNPs | ovary | testis | S4 |
| RNF20 | Z_random | | | 9 | SNPs, subtraction | both | | |
| ROBO1 | 1 | 107094521 | 107228509 | 19 | SNPs, coverage | ovary | | S5 |
| ROBO2 | 1 | 107529365 | 107979302 | 25 | SNPs | | | |
| RXRA | 17 | 8320685 | 8355067 | 14 | SNPs | | | S5 |
| SCRIB | 2 | 156239884 | 156325797 | 83 | SNPs | ovary | | S5 |
| SECISBP2L | 10 | 10159176 | 10193647 | 60 | SNPs, coverage | ovary | both | S5 |

| SHC4 | 10 | 10124441 | 10151124 | 11 | SNPs, coverage | | | S4 |
|------|-----|----------|----------|-----|-----------------|--------|--------|-----|
| SPHK1 | 18 | 7991834 | 7994408 | 2 | SNPs, coverage | testis | | |
| SRRT | Un | | | 16 | SNPs | both | | |
| SUGP2 | 28 | 4930094 | 4937971 | 33 | SNPs | ovary | both | S5 |
| SURF4 | 17 | 7682661 | 7693000 | 50 | SNPs | ovary | | S3 |
| TFEB | 26_random | 20475 | 21840 | 11 | SNPs | | | S5 |
| TIAM2 | 3 | 54800961 | 54890499 | NA | coverage | | | |
| TRIM71 | 2 | 60893878 | 60907039 | 159 | SNPs, subtraction | | | S1 |
| UBE2O | 18 | 7960889 | 7981633 | NA | coverage | | | |
| UGDH | 4 | 48113314 | 48126079 | 136 | SNPs, coverage, subtraction | ovary | ovary | S2 |
| UNC5C | 4 | 19035187 | 19126466 | 13 | SNPs, coverage | | | |
| Unnamed | Un | 124574513 | 124575553 | NA | coverage | | | |
| Unnamed | Un | 127129819 | 127130503 | NA | coverage | | | |
| Unnamed | 16_random | 26580 | 73126 | NA | coverage | | | |
| Unnamed | Un | 130103514 | 130104264 | NA | coverage | | | |
| Unnamed | Un | 50859565 | 50860210 | NA | coverage | | | |
| Unnamed | Un | 115355883 | 115358154 | NA | coverage | | | |
| Unnamed | Un | 124578595 | 124579326 | NA | coverage | | | |
| VEGFA | 3 | 31631385 | 31652650 | 34 | SNPs, coverage | both | | |
| WDR19 | 4 | 48204115 | 48240398 | 34 | SNPs | ovary | | S5 |
| ZWILCH | 10 | 19199771 | 19206407 | 8 | SNPs | ovary | | S5 |

Note: We were able to place only some genes on evolutionary strata due to our strict criteria for evaluating the maximum likelihood gene trees. The remaining genes lacked sequence information from several of the other sampled somatic genomes or had poorly resolved tree topologies.
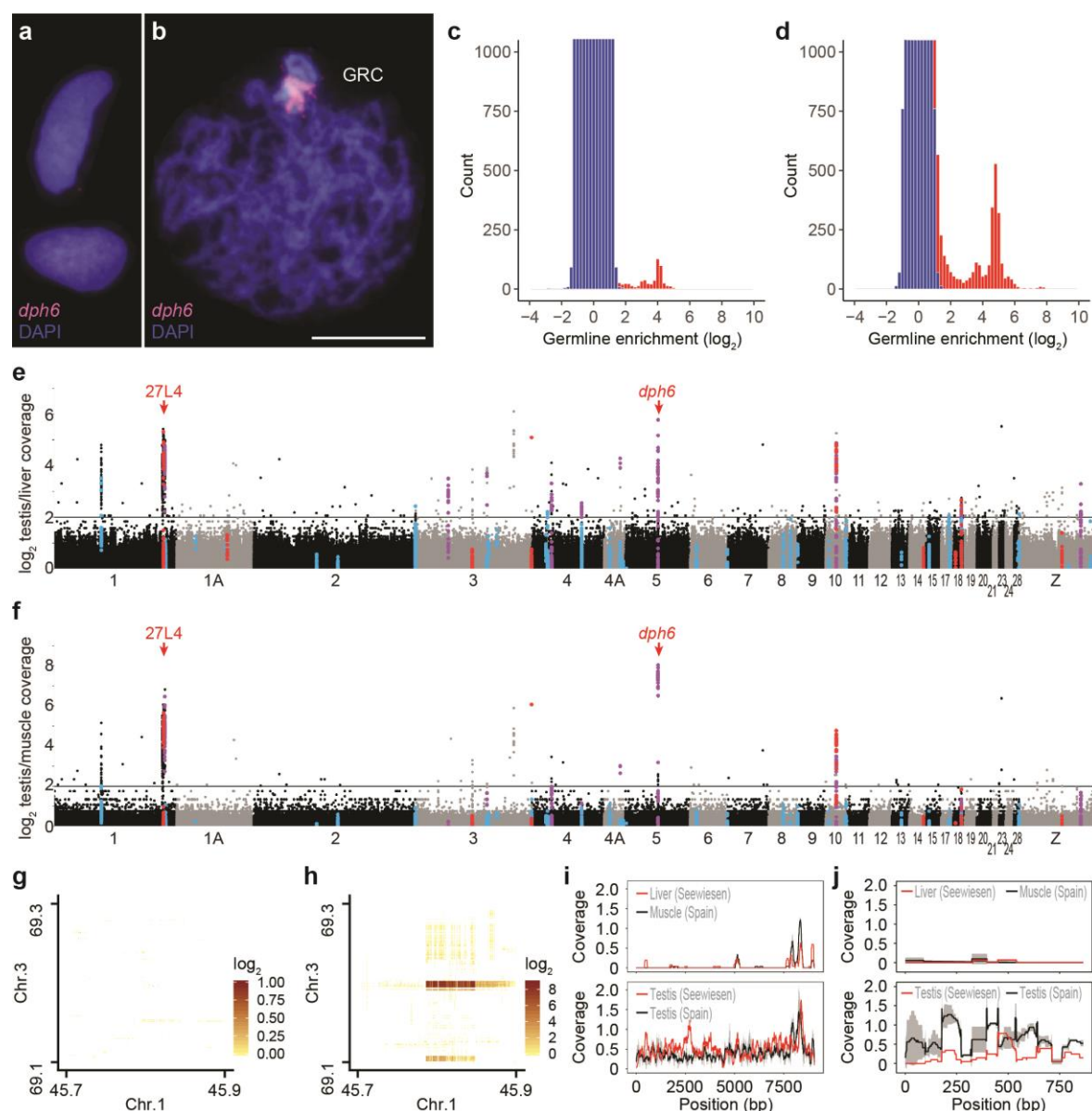
**Figure 1 | The zebra finch germline-restricted chromosome contains genes copied from many A chromosomes.**

**a-b**, Cytogenetic evidence for GRC absence in muscle (**a**) and GRC presence in the testis (**b**) of the same bird (Spain_1) using fluorescence *in-situ* hybridization (FISH) of our new GRC-ampliconic probe *dph6* (selected due its high germline/soma coverage ratio; *cf*. panels e-f). The scale bar indicates 10 μm. **c-d**, Comparison of germline/soma coverage ratios (red) for 1 kb windows with an expected symmetrical distribution (blue) indicates enrichment of single-copy regions in the germline, similar to lamprey[2] both in Spain (average of Spain_1 and Spain_2 coverage; PCR-free short reads) and Seewiesen (linked reads) samples. Y-axis is truncated for visualisation. **e-f**, Manhattan plot of germline/soma coverage ratios in 1 kb windows across chromosomes of the somatic reference genome taeGut2. Colours indicate high-confidence GRC-linked genes and their identification (red: coverage, blue: SNPs, purple: both; Table 1). Note that the similarities between Seewiesen (**e**) and Spain_1/Spain_2 averages (**f**) constitute independent biological replicates for GRC-ampliconic regions, as the data are based on different domesticated populations and different library preparation methods. Red arrows denote two FISH-verified GRC-amplified regions (*cf*. panel b)[8]. Only chromosomes >5 Mb are shown for clarity. **g-h**, Linked-read barcode interaction heatmaps of an inter-chromosomal rearrangement on the GRC absent in Seewiesen liver (**g**) but present in Seewiesen testis (**h**). **i-j**, Coverage plots of two examples of GRC-linked genes that are divergent from their A-chromosomal paralog, *trim71* (**i**) and *napa* (**j**)[9], and thus have very low coverage (normalized by total reads and genome size) in soma.
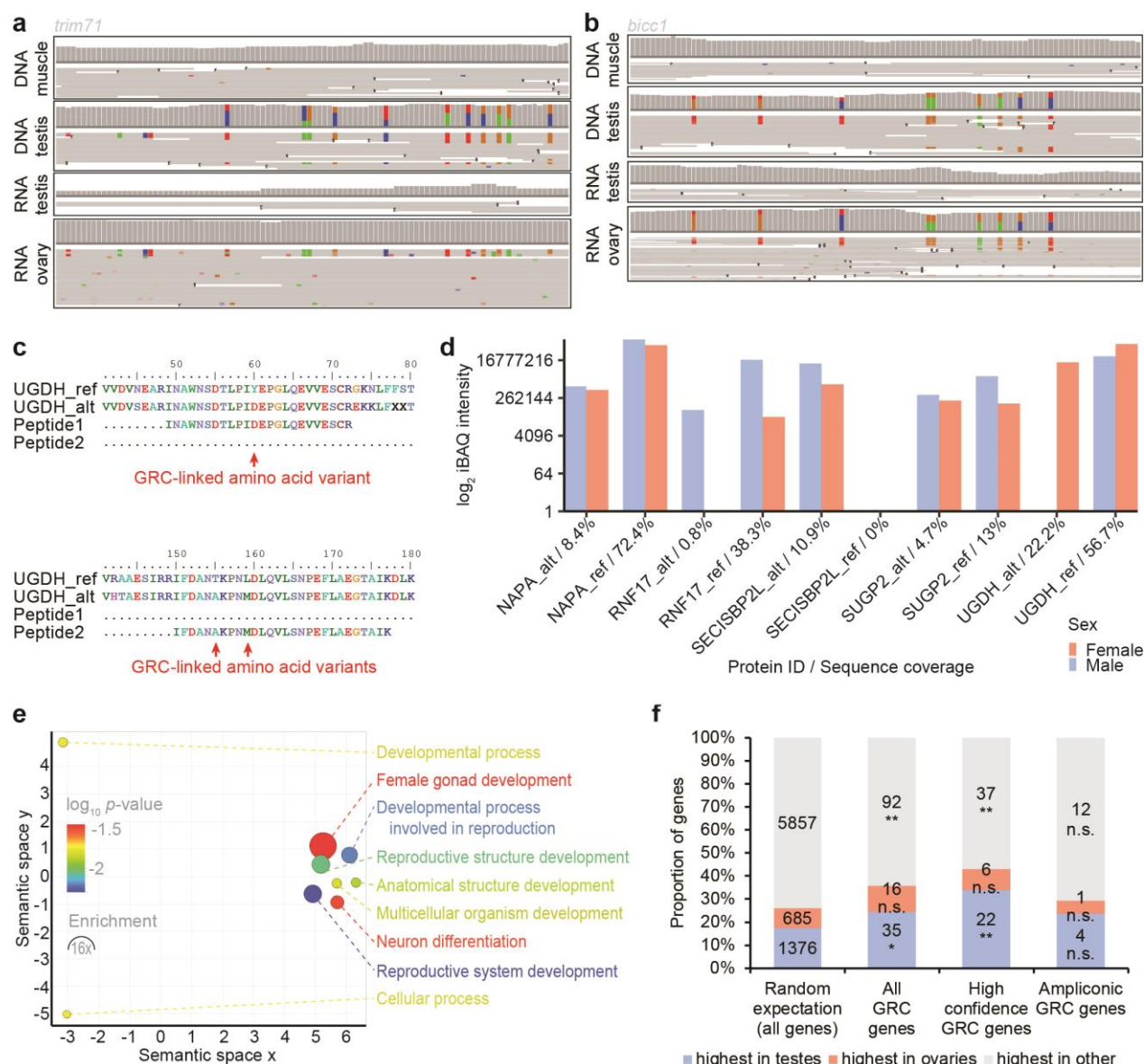
16

283
284 **Figure 2 | The zebra finch germline-restricted chromosome is expressed in male and female gonads. a-b**,
285 Comparison of coverage and read pileups for DNAseq from Spain_1 and Spain_2 testis/muscle, RNAseq data from
286 Spain_1 and Spain_2 testis, and available ovary RNAseq data[9]. Shown are 100-bp regions within *trim71* (**a**) and
287 *bicc1* (**b**). Colours indicate SNPs deviating from the reference genome taeGut2. **c**, Example alignment of proteomics
288 data showing peptide expression of the GRC-linked paralog of *ugdh* (alternative or 'alt'; *cf.* reference or 'ref'). **d**,
289 Proteomic evidence for GRC protein expression ('alt') in comparison to their A-chromosomal paralog ('ref') in 7
290 sampled testes and 2 sampled ovaries. Only GRC paralogs of RNF17 and UGDH appear to be expressed in a sex-
291 specific manner. **e**, Gene ontology term enrichment analysis of the 115 high-confidence GRC-linked genes (77
292 mapped gene symbols). Colours indicate the $\log_{10}$ of the false discovery rate *p*-value, circle sizes denote fold
293 enrichment above expected values. **f**, Expression evidence for orthologs of three different sets of GRC genes in testes,
294 ovaries, or other tissues in chicken[18]. Randomization tests show a significant enrichment for germline-expressed
295 genes among the 115 high-confidence GRC genes and all 267 GRC genes, but not the 38 ampliconic GRC genes.
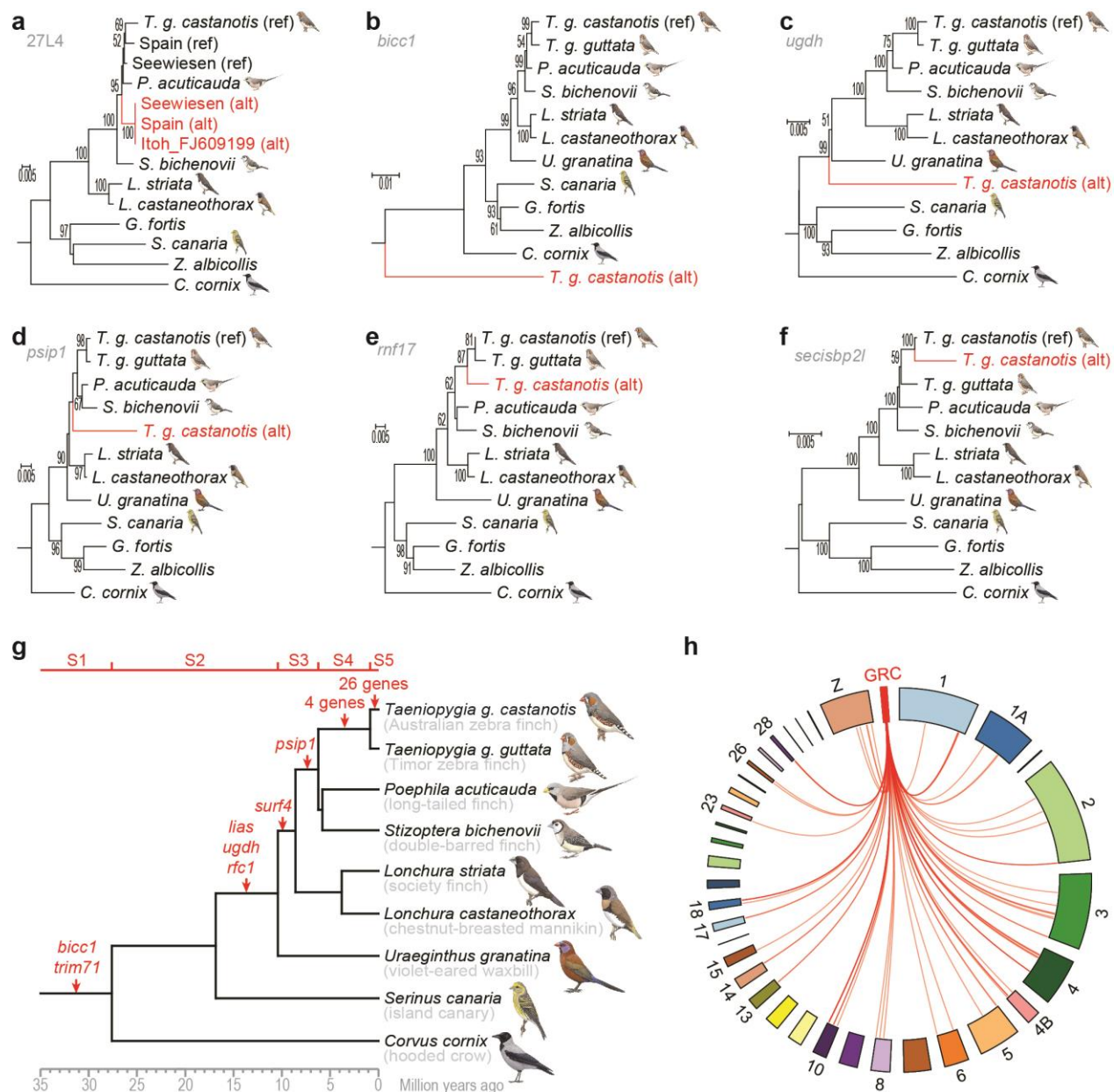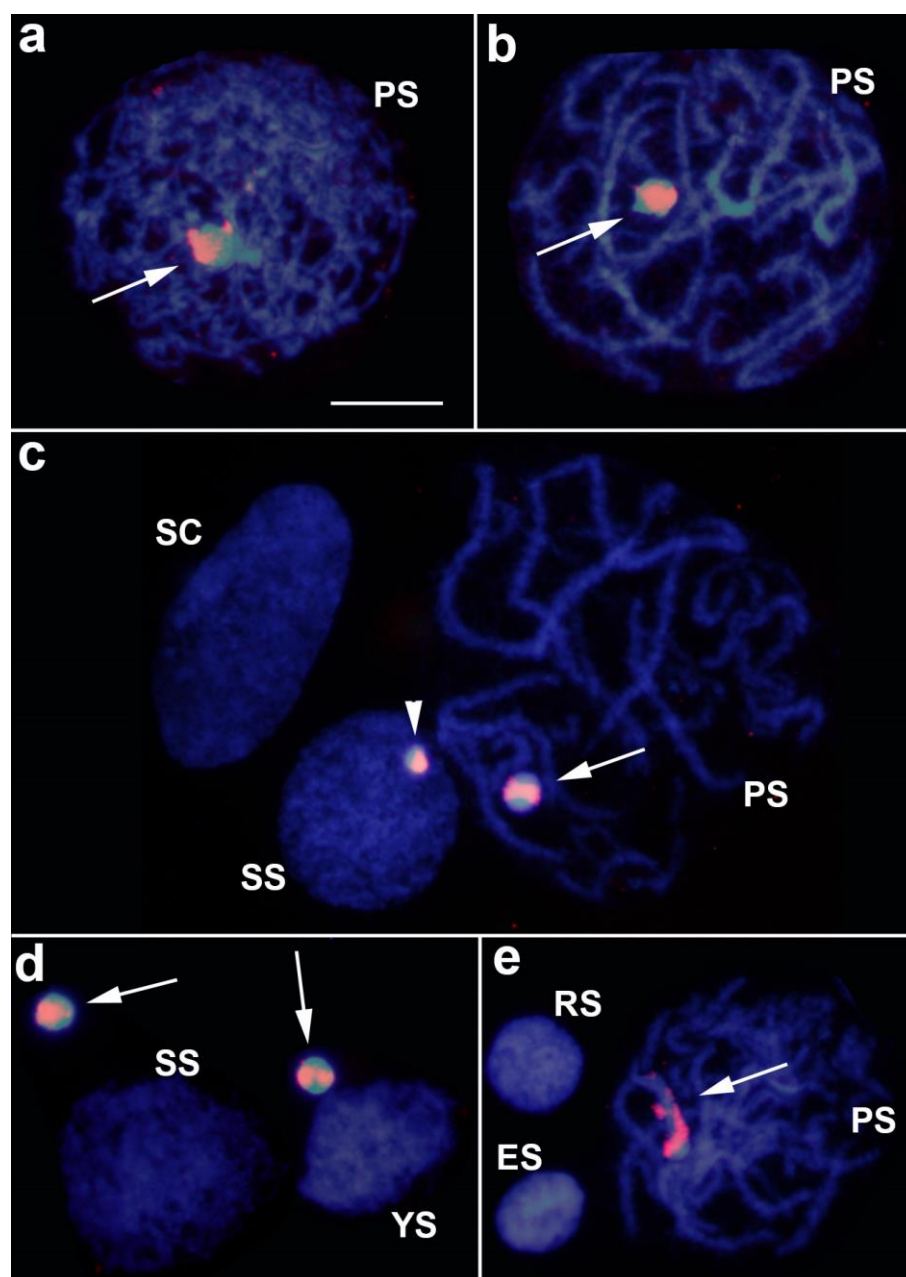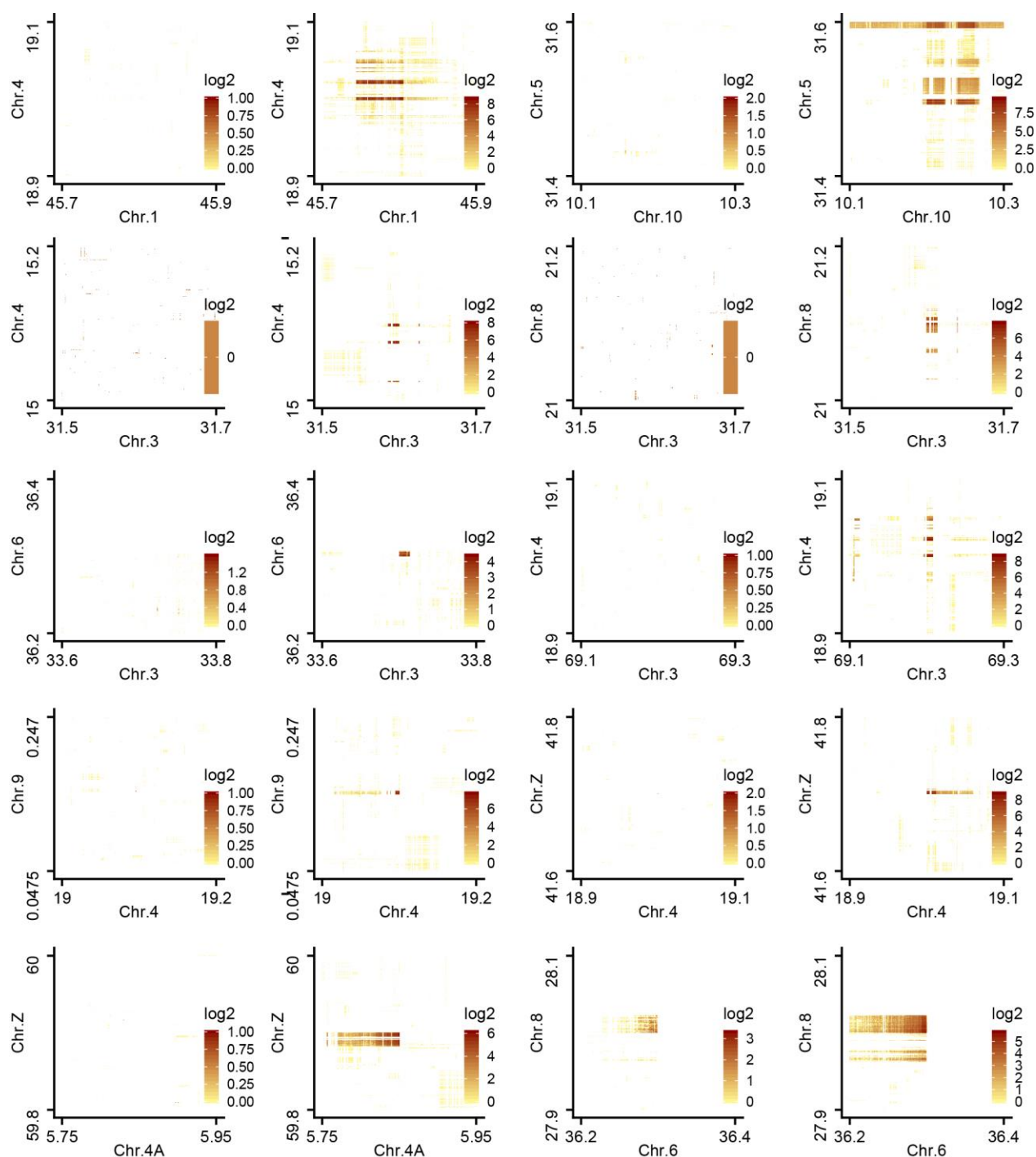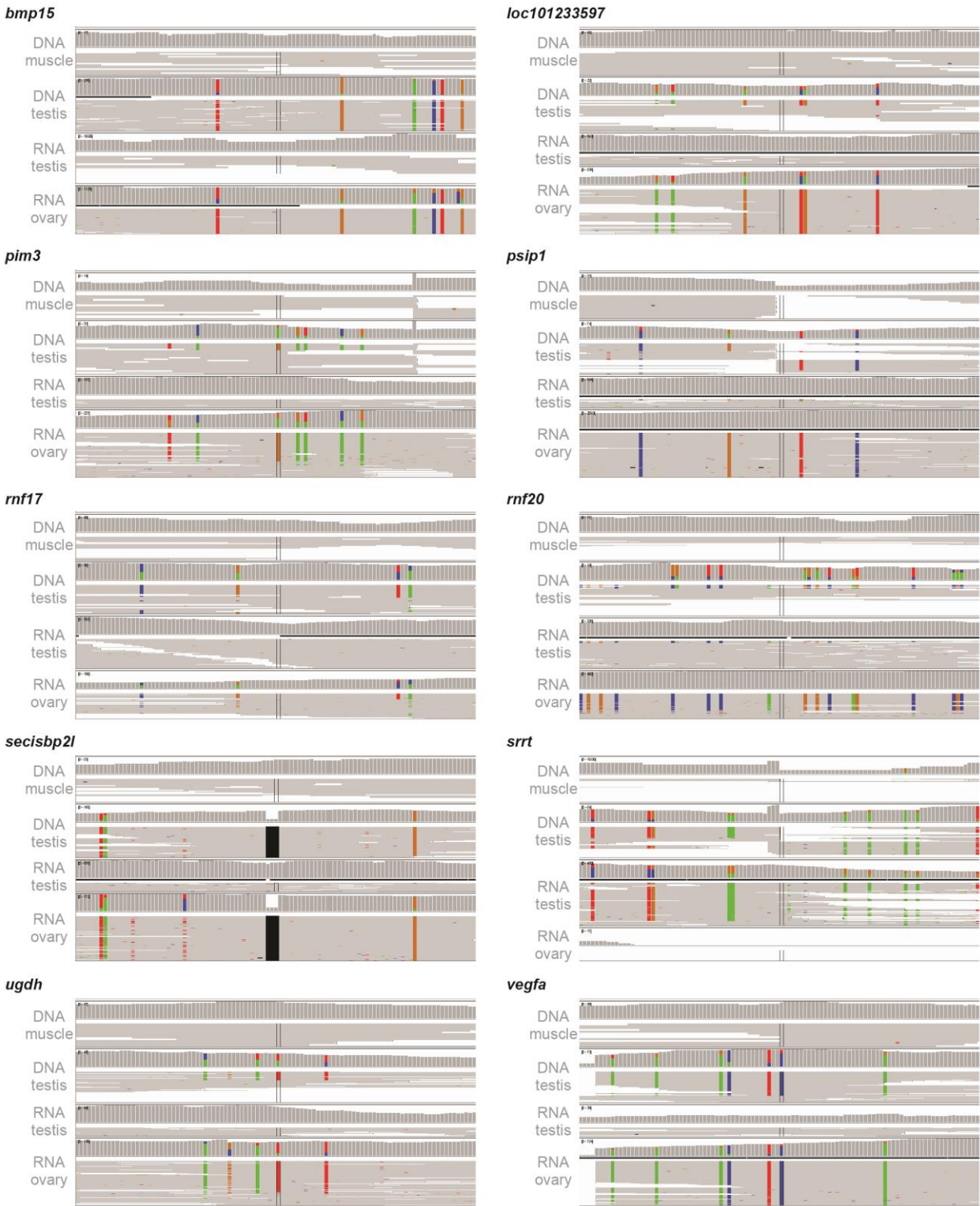
296

297

17

**Figure 3 | The zebra finch germline-restricted chromosome is ancient and highly dynamic. a**, Phylogeny of the intergenic 27L4 locus previously sequenced by Itoh et al.[8] suggests stable inheritance of the GRC paralog (alternative or 'alt' in red; *cf.* reference or 'ref') among the sampled zebra finches. **b-f**, Phylogenies of GRC-linked genes ('alt', in red; most selected from expressed genes) diverging from their A-chromosomal paralogs ('ref') before/during early songbird evolution (**b**; *bicc1*, stratum 1; *cf.* Extended Data Fig. 5), during songbird evolution (**c**; *ugdh*, stratum 2), during estrildid finch evolution (**d**; *psip1*, stratum 3), in the ancestor of the zebra finch species (**e**; *rnf17*, stratum 4), and in the Australian zebra finch subspecies (**f**; *secisbp2l*; stratum 5). The maximum likelihood phylogenies in panels a-f (only bootstrap values ≥50% shown) include available somatic genome data from estrildid finches and other songbirds. **g**, Species tree of selected songbirds showing the emergence of evolutionary strata (S1–S5) on the GRC (red gene names). Molecular dates are based on previous phylogenies[22,25]. Bird illustrations were used with permission from Lynx Edicions. **h**, Circos plot indicating A-chromosomal origin of high-confidence GRC-linked genes from 18 autosomes and the Z chromosome. Note that A-chromosomal paralogs of 37 genes remain unplaced on chromosomes in the current zebra finch reference genome taeGut2.
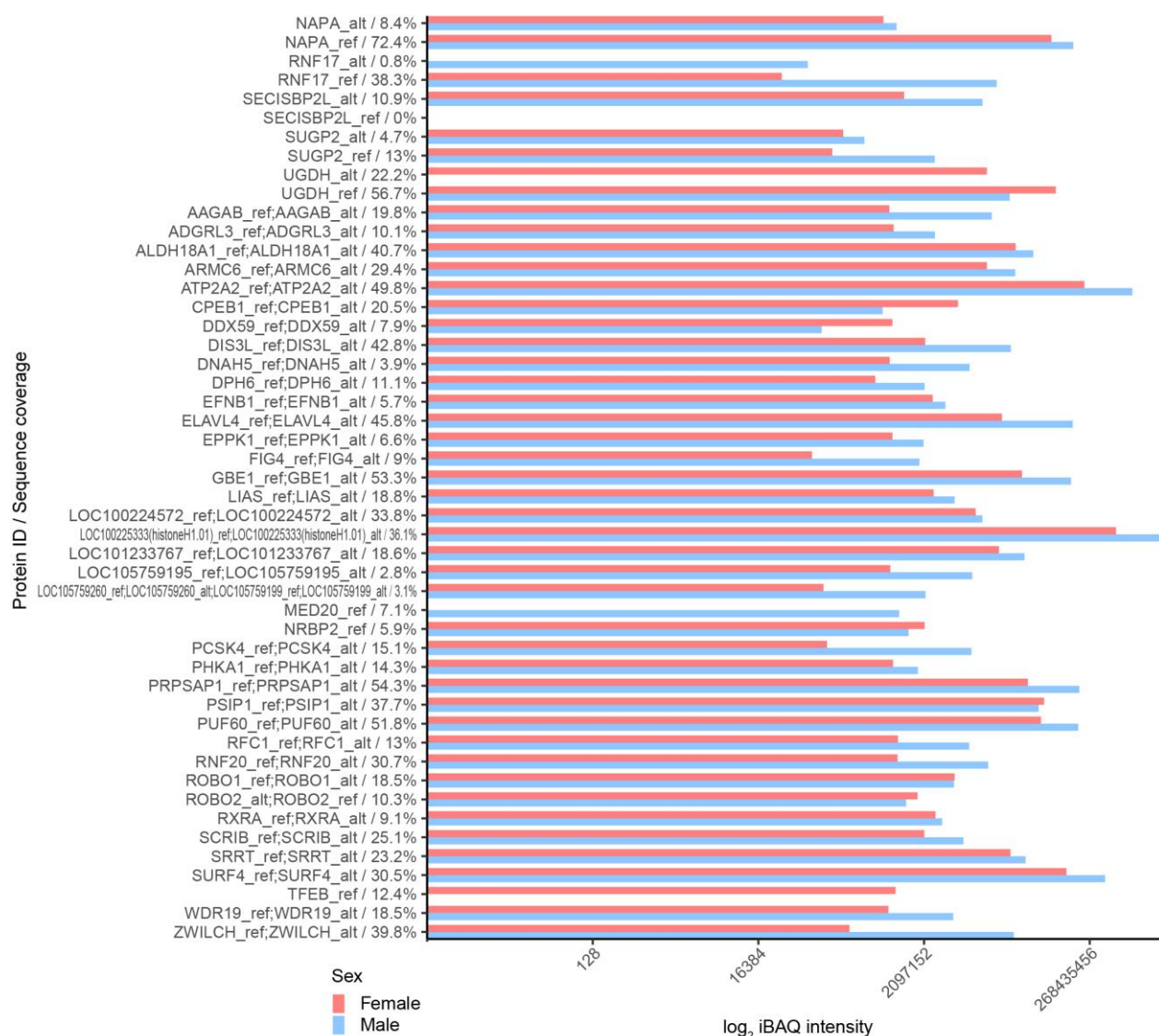
18

**Extended Data Figure 1 | FISH analysis in testis cells of the Spain_1 zebra finch individual using the *dph6* probe (red) counterstained with DAPI (blue).** Note the presence of primary (PS) and secondary (SS) spermatocytes, young spermatids (YS) and maturing spermatids at round (RS) and elongating (ES) stages. Also note that the *dph6* probe hybridizes with only part of the GRC chromosome (arrow), and this is apparent in PS at leptotene-zygotene (**a**), pachytene (**b-c**, **e**) and in GRCs which failed to integrate into the main nucleus of SS or YS cells (**d**), with no FISH signal in somatic cells (SC) indicating GRC absence in somatic structural testis cells (**c**). The half size of GRC in the SS cell in panel c, compared with that in the PS next to it and that those lying outside nuclei in panel d, indicates that GRC sometimes divides equationally in the first meiotic division (resulting in the half sized GRC body in panel c) but, in most cases, it divides reductionally yielding the large sized GRCs in panel d. Note that RS and ES nuclei in panel e lack FISH signal, indicating GRC absence. All photographs were made at the same magnification, and the scale bar in panel a indicates 10 μm.

19

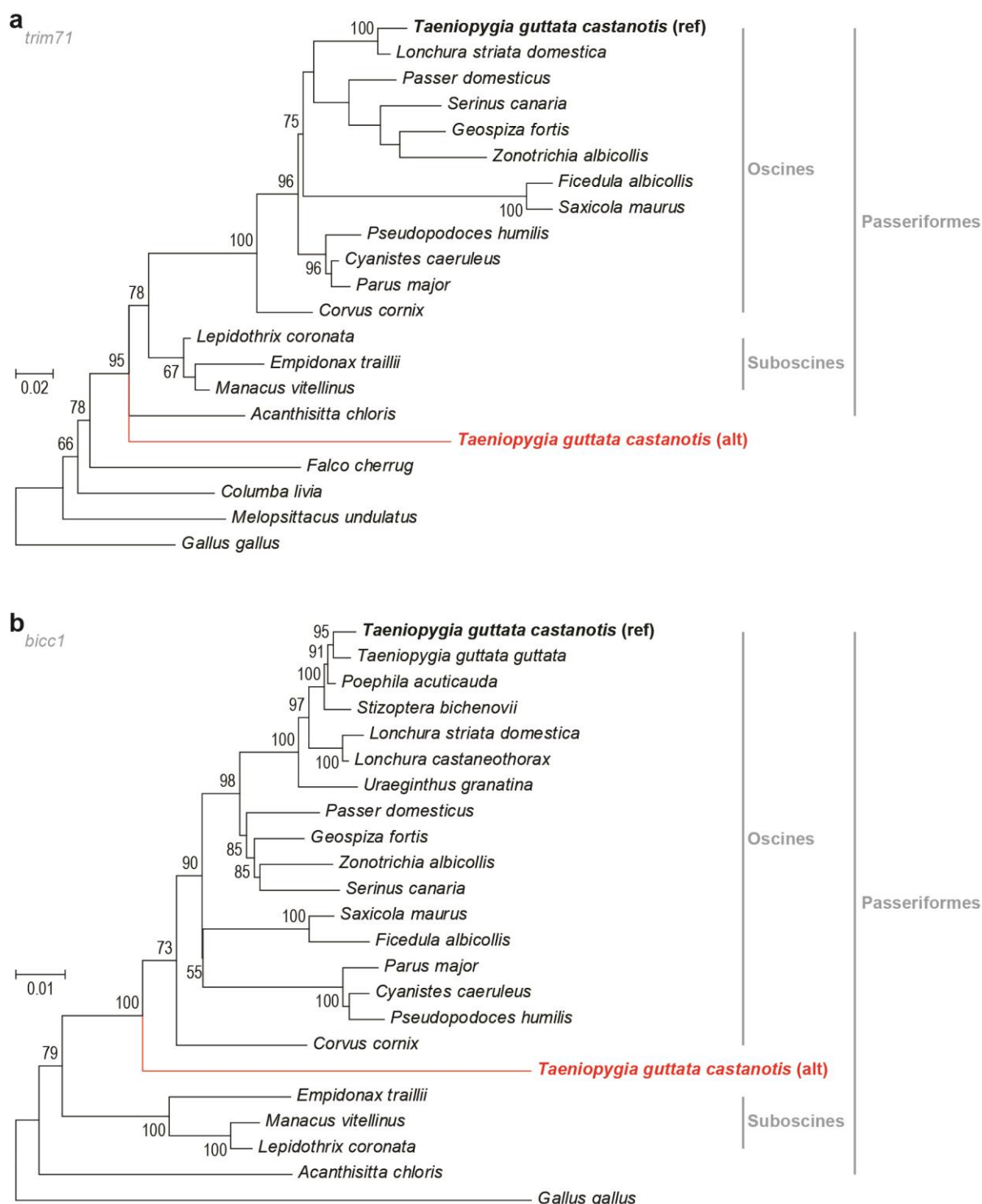325



326 **Extended Data Figure 2 | Testis-specific linked-read barcode sharing between A chromosomes indicates GRC**
327 **haplotypes.** Plots show side-by-side comparison of the inter-chromosomal barcode overlap for 200-kb regions for
328 the liver and testis, respectively (chromosome position scale in Mb). With the exception of the interaction between
329 chromosome 6 and chromosome 8 (bottom right) showing some background in the liver sample (potentially due to a
330 shared A-chromosomal rearrangement), all inter-chromosomal structural variants were testis-specific and thus
331 indicative of being on the same haplotype on the GRC. We exported barcode overlap matrices from the Loupe
332 browser for testis-specific structural variants called by LongRanger and plotted them in R (v. 3.5.1). We reassigned 0
333 values to "NA" (shown in white on the plot) and $\log_2$-transformed all values.

334

335

**Extended Data Figure 3 | Further examples for RNA expression of GRC-linked genes.** Comparison of coverage
and read pileups for DNAseq from Spain_1 and Spain_2 testis/muscle, RNAseq data from Spain_1 and Spain_2
testis, and available ovary RNAseq data[9]. Shown are 100-bp regions within 10 selected genes. Colours indicate SNPs
deviating from the zebra finch reference genome taeGut2.

340                                                                                                                21

341

**Extended Data Figure 4 | Proteomic evidence for functional GRC protein presence in zebra finch testes and ovaries.** The five proteins listed at the top are also shown in Fig. 2d. GRC paralogs are denoted by the 'alt' suffix, where A-chromosomal paralogs are denoted by the 'ref' suffix. Sequence coverage corresponds to the peptide coverage percentage of the reference protein sequence. Entries of only one protein identification have sufficient evidence at the peptide level to differentiate between the GRC and A-chromosomal paralogs due to coverage of non-identical regions between the both reference sequences; entries of more than one protein identification contain evidence of presence based solely on identical regions, thus cannot be differentiated at the proteomic level. Entries of only one protein identification without the corresponding 'alt' or 'ref' variant contain evidence that span the non-identical region only, thus the alternate variant need not be called.

351

22

**Extended Data Figure 5 | Gene trees of GRC-linked genes from stratum 1 and their A-chromosomal paralogs from broad taxon sampling imply GRC emergence in the ancestor of Passeriformes. a**, Maximum likelihood gene tree of *trim71* (partitioned for codon positions) suggesting GRC linkage in the ancestor of Passeriformes. **b**, Maximum likelihood gene tree of *bicc1* (only 3' UTR) suggesting GRC linkage in the ancestor of oscine songbirds.

23