Current Biology

Discovery of the First Germline-Restricted Gene by Subtractive Transcriptomic Analysis in the Zebra Finch, Taeniopygia guttata

Graphical Abstract



Authors

Michelle K. Biederman, Megan M. Nelson, Kathryn C. Asalone, Alyssa L. Pedersen, Colin J. Saldanha, John R. Bracht

Correspondence

jbracht@american.edu

In Brief

Biederman et al. report the first gene on the zebra finch germline-restricted chromosome (GRC): a gene encoding α -soluble NSF attachment protein (α -*SNAP*). Positive selection and higher expression in ovaries suggest a novel biological role, and the discovery of a somatic paralog in this first known instance of the gene occurring in multiple copies.

Highlights

- Discovery of the first germline-restricted gene, α-soluble NSF attachment protein
- Discovery of a somatic paralog (somatolog) of α-SNAP
- Positive selection and long branch length across α-SNAPs suggest novel function
- The α -SNAP pair exhibits sex-dimorphic expression, with GRC greater in ovaries





Discovery of the First Germline-Restricted Gene by Subtractive Transcriptomic Analysis in the Zebra Finch, *Taeniopygia guttata*

Michelle K. Biederman,¹ Megan M. Nelson,¹ Kathryn C. Asalone,¹ Alyssa L. Pedersen,¹ Colin J. Saldanha,¹ and John R. Bracht^{1,2,*}

¹Department of Biology, American University, 4400 Massachusetts Avenue NW, Washington, DC 20016, USA ²Lead Contact

*Correspondence: jbracht@american.edu

https://doi.org/10.1016/j.cub.2018.03.067

SUMMARY

Developmentally programmed genome rearrangements are rare in vertebrates, but have been reported in scattered lineages including the bandicoot, hagfish, lamprey, and zebra finch (Taeniopygia guttata) [1]. In the finch, a well-studied animal model for neuroendocrinology and vocal learning [2], one such programmed genome rearrangement involves a germline-restricted chromosome, or GRC, which is found in germlines of both sexes but eliminated from mature sperm [3, 4]. Transmitted only through the oocyte, it displays uniparental female-driven inheritance, and early in embryonic development is apparently eliminated from all somatic tissue in both sexes [3, 4]. The GRC comprises the longest finch chromosome at over 120 million base pairs [3], and previously the only known GRC-derived sequence was repetitive and non-coding [5]. Because the zebra finch genome project was sourced from male muscle (somatic) tissue [6], the remaining genomic sequence and protein-coding content of the GRC remain unknown. Here we report the first protein-coding gene from the GRC: a member of the α -soluble *N*-ethylmaleimide sensitive fusion protein (NSF) attachment protein $(\alpha$ -SNAP) family hitherto missing from zebra finch gene annotations. In addition to the GRC-encoded α -SNAP, we find an additional paralogous α -SNAP residing in the somatic genome (a somatolog)making the zebra finch the first example in which α -SNAP is not a single-copy gene. We show divergent, sex-biased expression for the paralogs and also that positive selection is detectable across the bird α -SNAP lineage, including the GRC-encoded α -SNAP. This study presents the identification and evolutionary characterization of the first protein-coding GRC gene in any organism.

RESULTS AND DISCUSSION

To identify genes from the germline-restricted chromosome (GRC), we adopted a subtractive transcriptomic approach. We

sequenced RNA from germline tissue of male and female adult birds, obtaining 10 million read pairs for each, and performed de novo assembly; we then performed computational elimination of sequences matching the published somatic (muscle) genome sequence [6], its raw (Sanger) read data, and a brain (somatic) transcriptome [7] (Figure 1A) to identify potential germline-limited sequences. During the filtering process, we identified 936 proteins having strong (1e-20 or better) matches to either the Swiss-Prot database or Pfam-A (thus, strong candidates for bona fide new genes) that are nevertheless missing from the current finch gene annotation (version 3.2.4 [8]). These new genes help fill in several important gaps in finch biology. For example, we uncovered a member of the DNA methyltransferase 1 (Dnmt1) family [9], an H1x linker histone [10], and the zeta subunit of the vesicle coat complex (COPI) [11], a member of the core eukaryotic orthologous family KOG3343.

The subtractive genomic pipeline uncovered a single GRC gene, a member of the α -soluble *N*-ethylmaleimide sensitive fusion protein (NSF) attachment protein (α-SNAP) family (hereafter the "GRC α-SNAP") (Figure 1A). Although our initial assembly captured a relatively short portion of the SNAP coding sequence that appeared to be an alternatively spliced isoform (Figure 1B, ii), we were able to reconstruct the full α -SNAP coding sequence by de novo assembly of a 94-million-read finch testis RNA sequencing (RNA-seq) dataset [12]. This assembled male-derived sequence matched the SNAP portion of the ovarian contig but encompasses a full SNAP coding sequence (Figure 1B, compare ii to v). We confirmed both isoforms by cloning and sequencing (Figure 1B). qPCR from a tissue panel of genomic DNA detected this gene at a statistically significant level only in testis (Figures 1C and S1A). Although the GRC α -SNAP was not detected in ovary DNA, we detected robust expression as RNA from this tissue (Figures 1D and S1C).

In the process of filtering the transcriptome data, we discovered a second α -SNAP gene. This one was filtered out from the raw Sanger reads in Figure 1A; thus, it is present in the somatic genome but is not present in the Sanger assembly [6]. Given we cannot use the term "gametolog," which refers to an autosomal copy of a sex-linked gene [13], we coin "somatolog" in reference to a somatic copy of a germline-limited gene. We suggest that the somatolog α -SNAP (the finch NAPA) underwent an ancient duplication event (possibly at the genesis of the GRC itself) forming a germline-restricted copy (NAPAG), which has subsequently undergone significant evolutionary divergence.



Figure 1. Discovery of a Paralogous α-SNAP Gene Pair

(A) Subtractive transcriptomic analysis used in this study.

(B) Overview of sequence comparison between assembled GRC (green) and somatolog (orange) α -SNAP sequences along with confirmation by cloning. (C) Genomic DNA qPCR analysis confirming *GRC* α -SNAP is only detected in testis or ovary (germline) tissue (primers F2+R2; see B). Error bars represent SEM. Two-way ANOVA identified testis signal (****p < 0.0001) as highly statistically significant, with n = 3 individuals of each sex tested.

(D) qRT-PCR analysis of expression of *GRC* α -*SNAP* showing strong ovary expression. Statistical significance was calculated with Student's two-tailed t test. Testis expression was significantly different from liver expression when using either oligo dT (p < 0.001) or random hexamer (p < 0.01). Ovary expression was significantly different from liver expression when using oligo dT (p < 0.01). Error bars represent SD of three measurements.

(E) qRT-PCR analysis of somatolog α -SNAP showing strong testis expression. Statistical significance was calculated with Student's two-tailed t test. Testis expression was significantly different from liver expression when using either oligo dT (p < 0.05) or random hexamer (p < 0.01). Error bars represent SD of three measurements.

See also Figure S1.

Expression of this paralogous gene system is sex biased. qRT-PCR of ovary and testis RNA revealed expression of the GRC gene is predominantly ovarian (Figure 1D), although gel analysis post-quantitation showed low-level detection also in testis (Figure S1C), as confirmed by our assembly of the gene from testis RNA-seq data and RT-PCR clone confirmation from both testis and ovary (Figure 1B). The somatolog α -SNAP is expressed in germlines and soma of both sexes, although most strongly in testis (Figure 1E).

We find that bird α -SNAP (NAPA) is in a particularly difficult-toassemble genomic location, leading to annotation problems for this gene family. Although β -SNAP genes have been deposited in GenBank for ten bird species, only two of these species had full-length α -SNAP genes available, and alignment of the protein sequences with each other and the zebra finch somatolog showed apparent discrepancies (Figure S2). Canary (Serinus canarius), accession number XP_009098415.2, aligns with α -SNAP of other species but has 20 central amino acids that are completely divergent, whereas society finch (*Lonchura striata domestica*), accession number XP_021401324, displays 30 altered amino acids at its carboxyl terminus (Figure S2). Ground tit (*Pseudopodoces humilis*) has an α -SNAP protein sequence (XP_005534295) that displays multiple problems: a 25 amino acid stretch is missing, as are the terminal 70 amino acids (Figure S2).



As we note above, no α -SNAP locus is present in the zebra finch Sanger assembly [6]; however, a recent (unannotated) PacBio haploid assembly [14] represents the locus as two allelic contigs, MUGN01000184.1 (386 kb) and MUGN01000615.1 (348 kb). Comparison of these with the α -SNAP scaffolds in both canary (unplaced scaffold NW_007931326.1, 84 kb) and society finch (unplaced scaffold NW_018657153.1, 496 kb) reveals that in all three species the exons are conserved, embedded within largely non-conserved repetitive micro- and minisatellite DNA. In both canary and society finch, assembly gaps obscure portions of exonic α -SNAP sequence, but no gaps exist in the two zebra finch PacBio contigs. The lack of high-confidence full-length a-SNAP protein annotations from other avian species (including the zebra finch) suggests that this region is problematic across birds, apparently both for assembly and for annotation.

To address this problem, we performed *de novo* RNA-seq assembly. Given that the zebra finch somatolog α -SNAP was robustly expressed within liver and testis (Figures 2A and 2C), we identified high-quality, deeply sequenced liver RNA-seq datasets from society finch (SRR5223631), canary (SRR2915372), and ground tit (SRR768235). (As in most birds, germline RNA-

Figure 2. Avian dN/dS Analysis of *α-SNAP* Genes

(A) Bayesian tree of bird SNAP proteins and dN/dS analysis of corresponding coding sequences. Red boxes represent β -SNAP; blue dots represent α -SNAP proteins. Branch numbers indicate posterior probabilities; scale bar represents substitutions per site. Branch letters A–I correspond to Table 1 for ω ratio (dN/dS) estimation of selection pressure.

(B) Analysis of ω for branch G using aBSREL [15] showing two selective regimes. Positive selection on this branch was statistically significant (p = 0.0045; Table 1).

(C) Analysis of ω for branch A using aBSREL [15] showing two selective regimes, with positive selection affecting 25% of sites but at a lower overall level than branch G. See also Figure S3.

seg datasets are not available for these species.) After Trinity assembly, we were able to retrieve a single full-length α -SNAP coding sequence for each bird that corrected the issues noted above (Figure S2). The canary, society finch, and ground tit α -SNAP genes were all identical, or within allelic variation, to the GenBank versions; however, the problematic regions have been replaced with sequences that, when translated, produce proteins that align confidently to other bird α -SNAPs (Figure S2). The transcriptomes of canary and ground tit also yielded full-length β -SNAP genes identical to those already deposited in GenBank, suggesting that our de novo assembly method is accurate (*β-SNAP*

is a brain-enriched, though not brain-exclusive, gene [16]). We also assembled RNA-seq datasets for great tit (*Parus major*) and golden-collared manakin (*Manacus vitellinus*), but these assemblies failed to yield full-length *SNAP* genes. Therefore, with the high-confidence zebra finch, canary, society finch, and ground tit α -SNAP gene sequences, we performed evolutionary tree reconstruction and analysis.

We aligned 14 sequences derived from 9 bird species (society finch was omitted for reasons described below) and built both Bayesian (Figure 2A) and maximum-likelihood (Figure S3) trees using chicken β -SNAP as outgroup. Both trees present nearly identical topologies and recover the α - and β -SNAP proteins as separate highly supported clades (Figures 2A and S3). However, although β -SNAP proteins are extremely well conserved among passerines, the α -SNAP proteins are much more divergent, located on extended branches (Figures 2A and S3). Indeed, the β -SNAP amino acid sequences of the canary, ground tit, society finch, and zebra finch are 100% identical, whereas manakin has only a single amino acid substitution (species-specific synonymous polymorphisms are present in their mRNAs). This suggests that β -SNAP genes are under significant purifying selection across passerines. In contrast, α -SNAP genes are



Figure 3. Multi-species Bayesian Tree, Confirming that Both SNAP Genes in Finch Are from the α -SNAP Family Red boxes, β -SNAP proteins; blue dots, α -SNAP proteins. Branch numbers indicate posterior probabilities; scale bar represents substitutions per site. See also Figure S2.

widely divergent in passerines, with 23 amino acid substitutions between the somatolog and canary α -SNAP and a somatologto-ground tit divergence of 44 amino acids. Remarkably, the interparalog, intra-zebra finch divergence is greater than that between the somatolog and all other passerine α -SNAP proteins (46 amino acids, ignoring the 8 amino acid deletion; Figure S1E). This results in an extremely long branch rivaling the one rooting the entire passerine clade (compare branches A and G in Figure 2A).

Society finch produced an unanticipated complexity: its α -SNAP consistently and confidently groups with the zebra finch GRC α -SNAP rather than somatic proteins (Figure 3). Society finch is the only passerine beside the zebra finch confirmed to have a GRC [17], but we derived this α -SNAP coding sequence from *de novo* assembled female liver (somatic) RNA-seq data. Furthermore, the gene is highly similar to a previous annotated version based on a blood-sourced (somatic) genome assembly (XP_021401324) (Figure S2). Due to the uncertainty surrounding this particular sequence, and the possibility that the unusual phylogenetic grouping is due to a long-branch attraction artifact [18, 19], we excluded society finch both from the trees in Figures 2A and S3 but include it in Figure 3.

We analyzed the bird-only phylogenetic tree (Figure 2A) for evidence of positive selection by analyzing the ratio of nonsynonymous mutations (dN) relative to synonymous (dS) mutations. When the dN/dS ratio, ω , is equal to 1, it implies the sequence is evolving neutrally-suggesting a loss of function on a coding sequence. Purifying selection-the weeding out of deleterious mutations to retain function-is indicated by ω less than 1, whereas positive selection-the promotion of specific amino acid changes due to advantageous function-is indicated by ω greater than 1 [20]. Branch models estimate ω for a whole protein (averaged across all amino acid sites), whereas branch-site models allow ω to vary across the amino acid sites at a specific branch of a phylogenetic tree [21]. This is a more sensitive method, because positive selection may only affect a few amino acids in a protein transiently during evolution, whereas most of the sites remain under purifying selection and mask the positive signal [20].

Analyzing the tree in Figure 2A, we found all branches (A–I) were estimated to have ω between 0 and 1, suggesting purifying selection (Table 1, branch model). However, we observed significant variation in ω estimates along lineages, with branches A and B in particular being elevated ($\omega = 0.548$ and 0.827; Table 1). Of

| Table 1. Analysis of dN/dS Ratios | | | | | | | | | | | | |
|-----------------------------------|---------------|------------------------------------|-------------------------|---|--|--|---|------------------------------------|-------------------------------------|---|---------|---|
| | PAML Analysis | | | aBSREL Analysis | | | | | | | | |
| | Branch Model | Branch-Site Model | | | | | | Branch-Site Model | | | | |
| Branch | ω Branch | ω ₂ (Positive Sites) | $ω_1$ (Purifying Sites) | Fraction Sites under Positive Selection $(\omega > 1)$ | Fraction Sites under Purifying Selection $(\omega < 1)$ | LRT Statistic (2∆InL) ^a | Significant Positively Selected Amino Acids (Posterior Probability) | ω ₂ (Positive Sites) | ω ₁ (Purifying Sites) | Fraction of Sites under Positive Selection | p Value | Significant after Multiple Hypothesis Testing? |
| A | 0.548 | 1.73 | 0.029 | 0.363 | 0.628 | 0.904 | 1 M (0.997), 14 N (0.997), 36 R (0.966), 107 R (0.999), 167 E (0.998), 266 W (0.998) | 3.03 | 0.249 | 0.25 | 0.1356 | no |
| В | 0.827 | 70.6 | 0.043 | 0.021 | 0.971 | 0.854 | _ | - | - | _ | - | - |
| С | 0.173 | 3.27 | 0.041 | 0.057 | 0.933 | 1.154 | 130 E (0.988) | - | - | _ | - | - |
| D | 0.048**** | 13.28 | 0.043 | 0.01 | 0.979 | 3.78 | 41 A (0.984) | - | - | _ | - | - |
| E | 0.103 | 999.00 (dS =0) | 0.042 | 0.02 | 0.968 | 13.807** | - | 15.2 | 0 | 0.024 | 0.0293 | no |
| F | 0.097**** | 1.21 | 0.039 | 0.098 | 0.89 | 0.035 | _ | - | _ | _ | - | _ |
| G | 0.128** | 4.53 | 0.042 | 0.06 | 0.928 | 0.707 | 1 M (0.974), 43 C (0.983), 171 R (0.980), 182 V (0953) | 2000 | 0.147 | 0.058 | 0.0045 | yes |
| Н | 0.066 | 1 | 0.044 | 0 | 0.988 | 0 | - | - | - | - | - | - |
| I | 0.053**** | 4.91 | 0.041 | 0.045 | 0.942 | 5.298 | 97 R (0.978) | - | - | - | - | - |

Analysis of dN/dS ratios (ω) for branches of the tree in Figure 2A, with bolded values indicating positive selection ($\omega > 1$). *p < 0.05, **p < 0.01, ****p < 10⁻⁴, corrected for multiple hypothesis testing.^aLRT, likelihood ratio test statistic, used in χ^2 . Critical values are 3.84 (5%) and 6.63 (1%).

the nine branches tested, only four were statistically significantly under purifying selection (Table 1), suggesting that the remaining branches are potentially under either relaxed purifying selection or positive selection at some sites. For the GRC and somatolog α -SNAP, this may be attributed directly to their paralogy, because genome-wide studies of gene duplication report relaxed purifying selection on paralogs, at least initially [22, 23]. However, the relaxed selection pressure is usually evolutionarily brief, reverting to a strongly purifying regime for both paralogs [22]. The GRC-somatolog a-SNAP divergence appears to be ancient by two measures: large amino acid divergence resulting in long branch lengths already noted (Figure 2A), and the synonymous (silent) mutations accruing between the copies, with pairwise dS = 0.26 by PAML. Most duplicate genes are lost (non-functionalized or turned into a pseudogene) by the time dS reaches a few percent [22], so the fact that both zebra finch genes produced by the *a-SNAP* duplication have been retained may indicate evolution of new function by the GRC copy.

We hypothesized that if the long branch lengths reflect selection for novel function, the elevated branch ω values (Table 1) might reflect a mixture of positive and purifying selection acting at different sites in the protein. We therefore evaluated branchsite models, which detect different selection pressures at specific branches across sites in a protein [21, 24]. PAML analysis uncovered positive selection on branches A, B, C, D, E, F, G, and I, all of which had some proportion of sites under $\omega_2 > 1$ (Table 1). The background purifying selection (ω_1) was found to be extremely consistent across branches and to account for the majority of sites on all branches tested (Table 1). Specific positively selected amino acids were identified by an empirical Bayesian approach [24] at posterior probability 0.95 or better for branches A, C, D, G, and I (Table 1). Branch E stands out with 0.02 of sites at an estimated ω_2 of 999, which means no synonymous substitutions were observed (dS = 0). Although dN/dS cannot be taken as a real value, likelihood ratios can still be accurately calculated for this branch, yielding a highly significant p <0.01 for positive selection, and branch I was significant to p < 0.05 prior to multiple testing correction (Table 1). Branch A leading to the GRC α-SNAP yielded a relatively modest ω_2 of 1.73, but it had by far the most sites under positive selection (0.363), possibly explaining why the branch ω was elevated ($\omega = 0.548$).

To confirm these findings, we ran the adaptive branchsite random effects likelihood (aBSREL) algorithm, which is similar to PAML but builds the tree and estimates the model complexity directly from the input sequence alignment [15]. Branch G showed statistically significant strong positive selection ($\omega_2 = 2000$ at 5.8% of sites, p = 0.0045; Table 1; Figure 2B), whereas E was also statistically significant (p = 0.03) before multiple hypothesis correction. (In all cases, we performed simple Bonferroni correction, which has been advocated in branch-site analysis [25], but may be too stringent [26, 27]; correction of Bayesian posterior probabilities is not required [28, 29].) We conclude that the positive selection along the phylogeny in Figure 2A is of extremely variable strength and distribution among sites. Branches G and E exhibit strong selection at 2%-6% of sites, whereas branch A (leading to the GRC α-SNAP) evidences a weaker positive selection across 25%-36% of sites (Table 1; Figures 2B and 2C).

To evaluate the wider evolutionary context of α - and β -SNAP genes, we aligned 16 α -SNAP and 20 β -SNAP proteins from

birds, reptiles, mammals, and fish. Consistent with the bird tree (Figure 2A), we recover α - and β -SNAP proteins as separate clades, and β -SNAPs have generally shorter branch lengths (Figure 3). Long β -SNAP branches occur in fish, specifically Atlantic herring (*Clupea harengus*) and great blue-spotted mudskipper (*Boleophthalmus pectinirostris*), which also display β -SNAP paralogy, the only cases outside the zebra finch α -SNAP duplication described in this work (Figure 3) in which the SNAP genes are duplicated.

The placement of fish and bird α -SNAP as sister clades is surprising (Figure 3). We do not have high confidence in this arrangement, as the branch support is lower. Instead, we attribute this grouping to the extremely long branch length of bird α -SNAP proteins creating long-branch attraction [18] and causing them to root basal to the mammal-reptile-chicken clade (Figure 3). This has been reported to be a risk of Bayesian reconstruction specifically in cases of rapidly evolving lineages with rate heterogeneity among sites [19]. However, a maximum-likelihood (RAxML) tree built from the same data displayed the same topology (not shown). Finally, we note the placement of society finch α -SNAP with the zebra finch GRC as a sister clade, an arrangement discussed above and which is extremely well supported, and may also be due to long-branch attraction.

In this work, we have identified the first gene from the GRC in the zebra finch, and the first case of α -SNAP paralogy in any organism. We confirmed this by searching the avian α -SNAP genes deposited in GenBank, representing 25 bird species, and noting that any duplicates we found were redundant copies of the same α -SNAP gene. To uncover potentially missed SNAP genes, we performed a tblastn search of the RefSeq passerine genomes and only uncovered single-copy α - and β -SNAP genes, consistent with the literature [16, 30].

The GRC-to-somatolog α -SNAP amino acid divergence (81% identity, 88% similarity) is comparable in scale to the divergence between zebra finch α - (somatolog) and β -SNAPs (72% identity, 90% similarity). This is reflected also by the branch leading to the GRC α -SNAP being nearly as long as the branches separating α - and β -SNAP clades (Figure 2A). Therefore, we cannot exclude the possibility that the GRC-encoded gene is a pioneering member of a new SNAP family [31]. Demonstrating this will require isolating more GRC *SNAP* genes from other birds, and to date germline genomic data are sorely lacking.

Because the duplicated gene in the zebra finch, the GRC α -SNAP, is present on a germline-limited sequence, the paralogy causes an effective doubling of the α -SNAP copy number in the germline only. Perhaps in response to this, the two paralogous genes have diverged to a high degree under positive selection. We also demonstrate that the two genes have sex-dimorphic expression in the germline, with the GRC α -SNAP more highly expressed in ovary than in testis. These data suggest the finch GRC is most likely playing an important biological role, in agreement with other studies showing that germline-restricted sequences are often involved in sex determination or germline function [1, 32], and we predict that more GRC-encoded genes are awaiting discovery. Finally, we note that if the gene duplication event leading to the zebra finch α -SNAP paralogy was the genesis of the GRC itself, our data imply that the GRC is relatively old and may be present in more bird lineages than originally expected.

STAR*METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- CONTACT FOR REAGENT AND RESOURCE SHARING
- EXPERIMENTAL MODEL AND SUBJECT DETAILS
- METHOD DETAILS
 - RNA extraction & sequencing
 - Sequencing
 - Error Correction and Assembly
 - Assembly of publically available RNA-seq datasets
 - dN/dS analysis
 - Phylogenetic Trees
 - Subtractive genomics for zebra finch
 - o qPCR and PCR
 - Reverse Transcription
- QUANTIFICATION AND STATISTICAL ANALYSIS
- DATA AND SOFTWARE AVAILABILITY

SUPPLEMENTAL INFORMATION

Supplemental Information includes three figures and can be found with this article online at https://doi.org/10.1016/j.cub.2018.03.067.

ACKNOWLEDGMENTS

The authors wish to acknowledge Naden Krogan for insightful comments on the manuscript, and Evgeny Bisk for Zorro computing cluster system administration. We are grateful to the Balakrishnan laboratory for the generation of two key transcriptomic resources used in this study: the testis (SRR2299402-4) and auditory lobule (SRS576610,11,12) RNA-seq datasets. The American University High Performance Computing System was funded in part by a grant from the National Science Foundation (BCS-1039497). This work was supported by NIH grant 1K22CA184297 (to J.R.B.), NIH NS 047267 (to C.J.S.), and an AU Faculty Research Support grant (to C.J.S. and J.R.B.).

AUTHOR CONTRIBUTIONS

M.K.B. performed tissue panel genomic DNA qPCR and qRT-PCR. M.M.N. performed RNA isolation, sequencing, error correction, assembly, and initial validation qPCR. K.C.A. performed codeml analysis on phylogenetic trees and made Figure 2A. A.L.P. and C.J.S. performed tissue isolation and initial RNA extractions. J.R.B. and C.J.S. conceived the study. C.J.S. provided biological material for sequencing. J.R.B. oversaw the study, performed assembly of public transcriptomic data and phylogenetic analysis, made most figures, and wrote the paper.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: November 8, 2017 Revised: February 12, 2018 Accepted: March 28, 2018 Published: May 3, 2018

REFERENCES

- Wang, J., and Davis, R.E. (2014). Programmed DNA elimination in multicellular organisms. Curr. Opin. Genet. Dev. 27, 26–34.
- Gurney, M.E., and Konishi, M. (1980). Hormone-induced sexual differentiation of brain and behavior in zebra finches. Science 208, 1380–1383.

- Pigozzi, M.I., and Solari, A.J. (1998). Germ cell restriction and regular transmission of an accessory chromosome that mimics a sex body in the zebra finch. Taeniopygia guttata. Chromosome Res. 6, 105–113.
- Pigozzi, M.I., and Solari, A.J. (2005). The germ-line-restricted chromosome in the zebra finch: recombination in females and elimination in males. Chromosoma 114, 403–409.
- Itoh, Y., Kampf, K., Pigozzi, M.I., and Arnold, A.P. (2009). Molecular cloning and characterization of the germline-restricted chromosome sequence in the zebra finch. Chromosoma 118, 527–536.
- Warren, W.C., Clayton, D.F., Ellegren, H., Arnold, A.P., Hillier, L.W., Künstner, A., Searle, S., White, S., Vilella, A.J., Fairley, S., et al. (2010). The genome of a songbird. Nature 464, 757–762.
- Balakrishnan, C.N., Lin, Y.C., London, S.E., and Clayton, D.F. (2012). RNAseq transcriptome analysis of male and female zebra finch cell lines. Genomics 100, 363–369.
- Washington University Genome Sequencing Center (2016). The zebra finch genome. https://www.ncbi.nlm.nih.gov/genome?term=taeniopygia %20guttata.
- Cheng, X., and Blumenthal, R.M. (2008). Mammalian DNA methyltransferases: a structural perspective. Structure 16, 341–350.
- Hergeth, S.P., and Schneider, R. (2015). The H1 linker histones: multifunctional proteins beyond the nucleosomal core particle. EMBO Rep. 16, 1439–1453.
- Futatsumori, M., Kasai, K., Takatsu, H., Shin, H.W., and Nakayama, K. (2000). Identification and characterization of novel isoforms of COP I subunits. J. Biochem. 128, 793–801.
- Singhal, S., Leffler, E.M., Sannareddy, K., Turner, I., Venn, O., Hooper, D.M., Strand, A.I., Li, Q., Raney, B., Balakrishnan, C.N., et al. (2015). Stable recombination hotspots in birds. Science 350, 928–932.
- García-Moreno, J., and Mindell, D.P. (2000). Rooting a phylogeny with homologous genes on opposite sex chromosomes (gametologs): a case study using avian CHD. Mol. Biol. Evol. 17, 1826–1832.
- 14. Korlach, J., Gedman, G., Kingan, S.B., Chin, C.S., Howard, J.T., Audet, J.N., Cantin, L., and Jarvis, E.D. (2017). De novo PacBio long-read and phased avian genome assemblies correct and add to reference genes generated with intermediate and short reads. Gigascience 6, 1–16.
- Smith, M.D., Wertheim, J.O., Weaver, S., Murrell, B., Scheffler, K., and Kosakovsky Pond, S.L. (2015). Less is more: an adaptive branch-site random effects model for efficient detection of episodic diversifying selection. Mol. Biol. Evol. 32, 1342–1353.
- Whiteheart, S.W., Griff, I.C., Brunner, M., Clary, D.O., Mayer, T., Buhrow, S.A., and Rothman, J.E. (1993). SNAP family of NSF attachment proteins includes a brain-specific isoform. Nature 362, 353–355.
- del Priore, L., and Pigozzi, M.I. (2014). Histone modifications related to chromosome silencing and elimination during male meiosis in Bengalese finch. Chromosoma 123, 293–302.
- Bergsten, J. (2005). A review of long-branch attraction. Cladistics 21, 163–193.
- Kolaczkowski, B., and Thornton, J.W. (2009). Long-branch attraction bias and inconsistency in Bayesian phylogenetics. PLoS ONE 4, e7891.
- Yang, Z. (2002). Inference of selection from multiple species alignments. Curr. Opin. Genet. Dev. 12, 688–694.
- Zhang, J., Nielsen, R., and Yang, Z. (2005). Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. Mol. Biol. Evol. 22, 2472–2479.
- Lynch, M., and Conery, J.S. (2000). The evolutionary fate and consequences of duplicate genes. Science 290, 1151–1155.
- Kondrashov, F.A., Rogozin, I.B., Wolf, Y.I., and Koonin, E.V. (2002). Selection in the evolution of gene duplications. Genome Biol. 3, RESEARCH0008.
- Yang, Z., Wong, W.S.W., and Nielsen, R. (2005). Bayes empirical Bayes inference of amino acid sites under positive selection. Mol. Biol. Evol. 22, 1107–1118.

- Anisimova, M., and Yang, Z. (2007). Multiple hypothesis testing to detect lineages under positive selection that affects only a few sites. Mol. Biol. Evol. 24, 1219–1228.
- Perneger, T.V. (1998). What's wrong with Bonferroni adjustments. BMJ 316, 1236–1238.
- Noble, W.S. (2009). How does multiple testing correction work? Nat. Biotechnol. 27, 1135–1137.
- Gelman, A., Hill, J., and Yajima, M. (2012). Why we (usually) don't have to worry about multiple comparisons. J. Res. Educ. Eff. 5, 189–211.
- Westfall, P.H., Johnson, W.O., and Utts, J.M. (1997). A Bayesian perspective on the Bonferroni adjustment. Biometrika 84, 419–427.
- Stenbeck, G. (1998). Soluble NSF-attachment proteins. Int. J. Biochem. Cell Biol. 30, 573–577.
- Clary, D.O., Griff, I.C., and Rothman, J.E. (1990). SNAPs, a family of NSF attachment proteins involved in intracellular membrane fusion in animals and yeast. Cell 61, 709–721.
- 32. Kloc, M., and Zagrodzinska, B. (2001). Chromatin elimination—an oddity or a common mechanism in differentiation and development? Differentiation 68, 84–91.
- 33. Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L., Raychowdhury, R., Zeng, Q., et al. (2011). Full-length transcriptome assembly from RNA-seq data without a reference genome. Nat. Biotechnol. 29, 644–652.
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local alignment search tool. J. Mol. Biol. 215, 403–410.

- Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics 25, 1754–1760.
- Yang, Z. (1997). PAML: a program package for phylogenetic analysis by maximum likelihood. Comput. Appl. Biosci. 13, 555–556.
- Yang, Z. (2007). PAML 4: phylogenetic analysis by maximum likelihood. Mol. Biol. Evol. 24, 1586–1591.
- Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinformatics 30, 1312–1313.
- Eddy, S.R. (2011). Accelerated profile HMM searches. PLoS Comput. Biol. 7, e1002195.
- Zhang, J., Kobert, K., Flouri, T., and Stamatakis, A. (2014). PEAR: a fast and accurate Illumina Paired-End reAd mergeR. Bioinformatics 30, 614–620.
- Yang, X., Dorman, K.S., and Aluru, S. (2010). Reptile: representative tiling for short read error correction. Bioinformatics 26, 2526–2533.
- 42. Yang, Z. (1998). Likelihood ratio tests for detecting positive selection and application to primate lysozyme evolution. Mol. Biol. Evol. *15*, 568–573.
- Katoh, K., Misawa, K., Kuma, K., and Miyata, T. (2002). MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. Nucleic Acids Res. 30, 3059–3066.
- Huelsenbeck, J.P., and Ronquist, F. (2001). MRBAYES: Bayesian inference of phylogenetic trees. Bioinformatics 17, 754–755.

STAR*METHODS

KEY RESOURCES TABLE

| REAGENT or RESOURCE | SOURCE | IDENTIFIER |
|---|------------------------------|--|
| Critical Commercial Assays | | |
| RNeasy Mini Kit | QIAGEN | Cat # 74104 |
| Nucleospin Tissue Kit | Machery-Nagel | Cat # 740952 |
| PowerSYBR qPCR mix | ThermoFisher | Cat # 4367659 |
| SuperScript III First Strand Synthesis System | ThermoFisher | Cat #18080051 |
| Deposited Data | | |
| Zebra finch testis RNA-seq data | [12] | SRA: SRR2299402 |
| Zebra finch testis RNA-seq data | [12] | SRA: SRR2299403 |
| Zebra finch testis RNA-seq data | [12] | SRA: SRR2299404 |
| Zebra finch auditory lobule RNA-seq data | [7] | SRA: SRS576610 |
| Zebra finch auditory lobule RNA-seq data | [7] | SRA: SRS576611 |
| Zebra finch auditory lobule RNA-seq data | [7] | SRA: SRS576612 |
| Society finch liver RNA-seq data | N/A | SRA: SRR5223631 |
| Canary liver RNA-seq data | N/A | SRA: SRR2915372 |
| Ground tit liver RNA-seg data | N/A | SRA: SRR768235 |
| Zebra finch ovary RNA-seg data | this paper | SRA: SRR6896649 |
| Zebra finch testis RNA-seg data | this paper | SRA: SRR6896648 |
| Zebra finch ovary RNA-seg data assembly | this paper | TSA: GGLD0000000 |
| 936 high-confidence zebra finch genes: nucleic | this paper | TSA: GGMT0000000 |
| acid / protein | | |
| Canary α-SNAP nucleotide / protein | this paper | GenBank: BK010484 |
| Society finch α -SNAP nucleotide / protein | this paper | GenBank: BK010485 |
| Ground tit α -SNAP nucleotide / protein | this paper | Genbank: BK010483 |
| Zebra finch GRC α -SNAP nucleotide / protein (NAPAG) | this paper | Genbank: MH263723 |
| Zebra finch somatolog α-SNAP nucleotide / protein (NAPA) | this paper | Genbank: MH263724 |
| Experimental Models: Organisms/Strains | | |
| Zebra Finch | Magnolia Bird Farm | N/A |
| Oligonucleotides | | |
| beta-actin_F 5'- TGGAGAAGAGCTACGAACTCCCTG - 3' | IDT | N/A |
| beta-actin_R 5'-GAAAGATGGCTGGAACAGGGCCTC - 3' | IDT | N/A |
| F1 5'-CGCGCCACCAAGCTCTTCAAGATG - 3' | IDT | N/A |
| R1 5'-CCTCGTGTTTGCTCTGCATCTGCAG - 3' | IDT | N/A |
| F2 5'-GGATGTGAGGGGCCGGAATTC - 3' | IDT | N/A |
| R2 5'-CCAACTCTCACCCTCCGGATC - 3' | IDT | N/A |
| F3 5'-GGCTCCTCGTTGCTGGAGGAG - 3' | IDT | N/A |
| R5 5'-GCCCTGGATGCTCCTCCTGATG- 3' | IDT | N/A |
| A 5'-GAGGCAGGCTCGGGGCTG - 3' | IDT | N/A |
| B 5'-TCCGCCTTTTTGAAGGCATTGCCAG - 3' | IDT | N/A |
| Software and Algorithms | | |
| Trinity (version 2.4.0) | [33] | https://github.com/trinityrnaseq/trinityrnaseq/ releases |
| Basic Local Alignment Search Tool (blast) | [34] | https://blast.ncbi.nlm.nih.gov/Blast.cgi? CMD=Web&PAGE_TYPE=BlastDocs&DOC_ TYPE=Download |
| TransDecoder (version 3.0.1) | Haas and Papanicolaou et al. | http://transdecoder.github.io |
| BWA (version 0.7.12) | [35] | https://sourceforge.net/projects/bio-bwa/files/ |

(Continued on next page)

| Continued | | | | | | |
|--------------------------|------------|--|--|--|--|--|
| REAGENT or RESOURCE | SOURCE | IDENTIFIER | | | | |
| PAML (version 4.9a) | [36, 37] | http://abacus.gene.ucl.ac.uk/software/ paml.html#download | | | | |
| aBSREL | [15] | http://datamonkey.org/absrel | | | | |
| RAxML | [38] | https://github.com/stamatak/standard-RAxML | | | | |
| Geneious (version 8.1.6) | Biomatters | http://www.biomatters.com | | | | |
| Hmmer (version 3.1b2) | [39] | http://hmmer.org/download.html | | | | |
| XLSTAT (version 2016) | Addinsoft | http://www.xlstat.com | | | | |

CONTACT FOR REAGENT AND RESOURCE SHARING

Further information and requests for resources and reagents should be directed to and will be fulfilled by the Lead Contact, John Bracht (jbracht@american.edu).

EXPERIMENTAL MODEL AND SUBJECT DETAILS

Adult zebra finches (*Taeniopygia guttata*) were obtained from a commercial breeder and housed in groups (15-25 per cage) in samesex aviaries. The colony room was maintained at 20C, 70% humidity and a 14:10 L:D cycle. Food, water and grit were available *ad libitum*. All animal husbandry was approved by the American University Animal Care and Use Committee.

RNA and DNA used in this study were extracted from germline and somatic tissue of five male and five female young adult birds.

METHOD DETAILS

RNA extraction & sequencing

Subjects were rapidly decapitated and tissue was removed and flash frozen on dry ice. Samples were then weighed and stored at -80 degrees until further processing. For RNA extraction, tissues were homogenized in 500 uL 100 mM Phosphate Buffer pH 7.4, and RNA was extracted from 100 uL of resultant homogenate using the RNeasy Mini Kit (QIAGEN) according to manufacturer instructions. The purity and concentration of each RNA sample was analyzed on a NanoDrop ND-100 spectrophotometer. Only extracts that exceeded a 260/280 ratio of 1.9 were used. Contaminating genomic DNA was eliminated by treatment with Turbo DNase (ThermoFisher Cat #AM2238) and 15-20 μ g of RNA was submitted to Eurofins Genomics (Huntsville AL). Total DNA was purified from homogenates using the Nucleospin Tissue kit (Macherey-Nagel, Düren, Germany) according to manufacturer's instructions.

Sequencing

Paired-end, expression-normalized, and strand-specific Illumina sequencing was performed by Eurofins Genomics (Huntsville AL). Read lengths were 300 basepair (bp) from a MiSeq, and the total number of read pairs obtained was 10,704,971 for Ovary and 9,703,220 for Testis.

Error Correction and Assembly

After sequencing the paired reads were stitched together with PEAR [40] to generate high-quality merged raw reads with a mean length of 300bp. Read error correction was performed using Reptile [41], followed by assembly on AU's Zorro High Performance Computing Cluster using Trinity [33] run in default mode and specifying the –SS_lib_type parameter for strand-specific libraries. Following assembly the longest open reading frames were identified using TransDecoder.LongOrfs.

Assembly of publically available RNA-seq datasets

Zebra finch testis: Datasets SRR2299402, SRR2299403, and SRR2299404 were downloaded from NCBI's Sequence Read Archive database (https://www.ncbi.nlm.nih.gov/sra). The fastq files were combined into a single file and Trinity was run using default settings and the '-trimmomatic' flag.

For society finch dataset SRR5223631 was downloaded (177 million reads) and assembled with the '-trimmomatic' and '-SS_lib_ type FR' flags.

For canary, dataset SRR2915372 was downloaded (123 million reads) and assembled with the '-trimmomatic' flag. For ground tit, dataset SRR768235 was downloaded (28 million reads), and assembled with the '-trimmomatic' flag.

dN/dS analysis

PAML Branch model

Codeml (PAML v.4.9) [36, 37] was used to estimate ω using the branch model setting (runmode = 0, seqtype = 1, model = 2, NSsites = 0). Branches to be estimated were specified in the newick tree file (for the Bayesian tree), one at a time. Each branch

was estimated twice: once with a neutral model (above settings plus fix_omega = 1 and omega = 1) and using a purifying selection model (fix_omega = 0, omega = 1). The p values were determined using the likelihood ratio test (LRT) statistic 2ΔI [42] compared against χ^2 with critical values of 3.84, 5% significance level, and 6.63, for 1% significance [31]. Correction for multiple hypothesis testing was performed.

PAML Branch-sites model

Branch-sites ω were estimated by adding NSsites = 2 to the Codeml control file and estimating one branch at a time. The p values were calculated as for the branch model by LRT statistic.

PAML Pairwise

A pairwise alignment of GRC and somatolog coding sequences was provided to Codeml with runmode = -2 and CodonFreq = 2. aBSREL

For aBSREL [15] the 14-sequence bird alignment that was used for tree building was input into the online interface (http://www. datamonkey.org). The relevant foreground branches were selected as indicated in Table 1.

Phylogenetic Trees

All sequences for comparison were obtained from NCBI or assembled de novo and curated for length, with short (incomplete) sequences discarded. Alignment was performed using the MAFFT algorithm [43] implemented within the Geneious software package (http://www.biomatters.com). The tree was generated with Mr. Bayes [44] in Geneious, using the Rate Matrix = equalin and Rate Variation = invgamma settings. The maximum likelihood (ML) tree was built using RAxML version 8.2.11 [38] with the GAMMA JTT protein model and 200 bootstrap replicates. The outgroup was chicken β -SNAP. Acccession numbers used are given below.

Accession numbers for α-SNAP mRNA are: human-NM 003827, mouse-NM 025898, rat-NM 080585, chicken-XM 015272486, painted turtle- XM_005310524, western clawed frog- NM_001011280, African clawed frog- NM_001092405, Atlantic herring-XM_012832758, Asian sea bass- XM_018665555, zebrafish- NM_199766, and great-blue spotted mudskipper- XM_020928551.

Accession numbers for α-SNAP proteins are: human- NP_003818.2, mouse- NP_080174.1, rat- NP_542152.1, chicken-XP_015127972.1, painted turtle- XP_005310581.1, western clawed frog- NP_001011280.1, African clawed frog- NP_001085874.1, Atlantic herring- XP 012688212.1, Asian sea bass- XP 018521071.1, zebrafish- NP 956060.1, and great-blue spotted mudskipper-XP 020784210.1.

Accession numbers for β-SNAP mRNA are: human- NM_001283018, mouse- NM_019632, rat- NM_001191966, chicken-NM_001199430, zebra finch- XM_002199762, ground tit- XM_005525483, canary- XM_009093739, rock dove- XM_005513170, downy woodpecker- XM_009902073, eagle- XM_010574905, Japanese quail- XM_015856683, golden-collared manakin-XM 018077509, western clawed frog- NM 001079098, zebrafish- NM 001080702, Atlantic herring- XM 012826735 and XM_012838056, African clawed frog- XM_018265067, and great-blue spotted mudskipper- XM_020921395 and XM_020933727.

Accession numbers for β-SNAP proteins are: human- NP_001269947.1, mouse- NP_062606, rat- NP_001178895.1, chicken-NP_001186359.1, zebra finch- XP_002199798.1, ground tit- XP_005525540.1, canary- XP_009091987.1, rock dove-XP_005513227.1, downy woodpecker- XP_009900375.1, eagle- XP_010573207.1, Japanese quail- XP_015712169.1, golden-collared manakin- XP_017932998.1, western clawed frog- NP_001072566.1, zebrafish- NP_001074171.2, Atlantic herring- XP_012682189.1 and XP_012693510.1, African clawed frog- XP_018120556.1, and great-blue spotted mudskipper- XP_020777054.1 and XP_020789386.1.

Subtractive genomics for zebra finch

Phase 1

The Sanger finch genome (GCA_000151805.1 Taeniopygia_guttata-3.2.4) and the mitochondrial sequence (MT) were downloaded from NCBI and combined into a single Basic Local Alignment Search Tool (blast) nucleotide database [34]. The Trinity ovary and testis assemblies were used as queries for local blastn against this combined genome+MT database, with default settings in order to provide maximal confidence in the remaining sequences' uniqueness. A custom python script was used to segregate the non-matching sequences (i.e., those with no blastn matches). Open reading frames were identified using TransDecoder.LongestOrfs,(supplied with Trinity software package) and we used custom python scripts to remove redundant protein isoforms by selecting for the longest protein-coding sequence from each gene. Potential protein-coding homologs were identified by 1) blastp against the uniprot-swissprot database (evalue 1e-20) or 2) Hmmer3.1b2 [39] search against the Pfam-A database (Pfam 28.0) (also requiring evalue 1e-20). Phase 2

The 936 proteins identified in Phase 1 were more stringently filtered by blast against the raw Sanger data, downloaded from the NCBI Trace Database (ftp://ftp-private.ncbi.nlm.nih.gov/pub/TraceDB/taeniopygia_guttata/). Default tblastn settings were used but we checked to confirm that evalues were highly significant and represented true matches. For example, we filtered out 520 Ovary hits (out of an initial set of 598, see Figure 1) giving 78 potential GRC genes. Of the 520 blast hits against Sanger raw reads, 517 (99.4%) occurred with an e-value of 2e-6 or better. Similarly our blast against Sanger raw reads filtered out 614 from a Testis set of 705 (keeping 91 genes) and of those blast hits 605 (98.5%) were of evalue 1.25e-6 or better.

This dataset was further filtered by mapping raw reads from a very large Auditory Lobule (brain) dataset generated by the Balakrishnan lab [7] (SRA archive SRS576610, SRS576611, and SRS576612), totaling approximately 70 million reads, onto the germline gene coding sequences with BWA [35] (bwamem, default settings). We eliminated any candidates with matching reads from the AL read mapping. Remarkably, this eliminated all but 8 of the 78 ovary transcripts: six were viral in nature (suggesting an unrecognized and apparently asymptomatic infection); of the remaining two, one was clearly repetitive in sequence and not considered further. The remaining gene was the novel SNAP protein (TR30145) confirmed to be GRC derived based on qPCR off genomic DNA and described further here. The testis dataset did not yield any GRC genes but yielded a contig encoding the somatolog α -SNAP, which was also assembled independently in the ovary transcriptome.

qPCR and PCR

Unless otherwise noted, all qPCR reactions (PowerSYBR, ThermoFisher Cat # 4367659) were run as a 2-stage cycle with 95° C for 10 min initial melt, then 40 cycles of 95° C for 30 s, 60° C for 1 min and measurement of DNA concentration. Primers F1 + R1 cannot be used for qPCR off genomic DNA owing to a 689bp intron situated between them, necessitating the construction of primers A + B used instead. To gain specificity with the A+B primer set required a customized 2-step cycle of 95° C for 30 s, followed by 64° C for 10 s (still run for 40 total cycles).

All qPCR signal was measured relative to actin by Δ Ct: we calculated average and standard deviation of $2^{-(gene\ Ct\ -\ \beta-actin\ Ct)}$ for all cases. Statistical significance was measured by Student's 2-tailed t test or 2-way ANOVA.

Normal (nonquantitative) PCR was carried out using AccuStart II polymerase (QuantaBio, Beverly, MA) and used according to manufacturer's instructions, with annealing at 58°C, extension for 1 min, and 35 cycles. Template was cDNA constructed as described below.

Reverse Transcription

SuperScript III First Strand Synthesis System (ThermoFisher Cat #18080051) was used in accordance with manufacturer's instructions, with 4 µg of total RNA that had been DNase-treated with Turbo DNase (ThermoFisher, Cat #AM2238) and phenol extracted. cDNAs were diluted 10x prior to use. Minus-RT controls were always tested in parallel to ensure no contaminating genomic DNA was present in the samples.

QUANTIFICATION AND STATISTICAL ANALYSIS

Statistical tests for Figure 1 are described in the legend of that figure and include Student's t test and ANOVA. For Figures 1C and S1A, three birds (n = 3) were tested per sex, and the SNAP / actin ratio was measured by qPCR for each tissue in triplicate, yielding nine overall measurements per tissue. The graph shows the average and standard error of the mean for these nine measurements. Two-way ANOVA was performed with the XLSTAT (http://www.xlstat.com) Excel add-on software package. For Figures 1D and 1E the graphs represent the average and error bars represent standard deviation of triplicate measurements, with statistical significance obtained by Student's 2-tailed t test.

For Figure S1 analysis was identical to Figure 1 except that for each of n = 3 birds per sex, the tissue was measured six times, yielding 18 overall measurements per tissue. The graph represents average and standard error of the mean for these 18 measurements. All other analysis as in Figure 1.

For Figure 2 the statistical significance was obtained by PAML and by aBSREL; however PAML required running twice per branch (once for null and alternative), obtaining likelihood ratios, and testing these ratios by chi-square as described in Method Details.

DATA AND SOFTWARE AVAILABILITY

The zebra finch ovary and testis RNA-seq reads have been deposited in SRA under accession number SRR6896649 and SRR6896648 and the assembled data in TSA under accession number GGLD00000000.

The 936 high-confidence genes identified in this study have been deposited in GenBank under accession number GGMT00000000.

The α -SNAP genes from zebra finch GRC α -SNAP, somatolog, canary, society finch, and ground tit have been deposited in GenBank under accession numbers MH263723, MH263724, BK010484, BK010485, and BK010483, respectively.