

# The effective population size modulates the strength of GC biased gene conversion in two passerines

Henry J Barton<sup>1,2,\*</sup> and Kai Zeng<sup>1,\*</sup>

<sup>1</sup>Department of Animal and Plant Sciences, The University of Sheffield, Western Bank, Sheffield, S10 2TN, UK

<sup>2</sup>Organismal and Evolutionary Biology Research Programme, Viikinkaari 9 (PL 56), University of Helsinki, Helsinki, FI-00014, Finland

\*Corresponding authors: [henry.juho.barton@gmail.com](mailto:henry.juho.barton@gmail.com), [kzeng@sheffield.ac.uk](mailto:kzeng@sheffield.ac.uk)

## Abstract

Understanding the determinants of genomic base composition is fundamental to understanding genome evolution. GC biased gene conversion (gBGC) is a key driving force behind genomic GC content, through the preferential incorporation of GC alleles over AT alleles during recombination, driving them towards fixation. The majority of work on gBGC has focussed on its role in coding regions, largely to address how it confounds estimates of selection. Non-coding regions have received less attention, particularly in regard to the interaction of gBGC and the effective population size ( $N_e$ ) within and between species. To address this, we investigate how the strength of gBGC ( $B = 4N_e b$ , where  $b$  is the conversion bias) varies within the non-coding genome of two wild passerines. We use a dataset of published high coverage genomes (10 great tits and 10 zebra finches) to estimate  $B$ , nucleotide diversity, changes in  $N_e$ , and crossover rates from linkage maps, in 1Mb homologous windows in each species. We demonstrate remarkable conservation of both  $B$  and crossover rate between species. We show that the mean strength of gBGC in the zebra finch is more than double that in the great tit, consistent with its twofold greater effective population size.  $B$  also correlates with both crossover rate and nucleotide diversity in each species. Finally, we estimate equilibrium GC content from both divergence and polymorphism data, which indicates that  $B$  has been increasing in both species, and provide support for population expansion explaining a large proportion of this increase in the zebra finch.

**Keywords**— GC biased gene conversion, gBGC, effective population size, equilibrium GC, population expansion

## 28 Significance statement

29 Understanding the forces that change the nucleotide base composition of genomes is central to understanding their  
30 evolution. One such force is GC biased gene conversion, a process that during recombination converts some heterozy-  
31 gous base positions to homozygous. This process is more likely to convert adenine and thymine bases to guanine  
32 and cytosine bases than the other way around, hence is GC biased. This increases the frequency of GC alleles in a  
33 way similar to positive selection. This process has largely been studied within protein coding regions, and not often  
34 compared between species. We measure its strength in the non-coding areas of the genomes of two bird species,  
35 showing it to be stronger in the species with the larger population size.

## 36 Introduction

37 A large proportion of many organisms' genomes are non-coding; 99% in humans, 80% in *Drosophila melanogaster*,  
38 73% in *Caenorhabditis elegans* and 71% in *Arabidopsis thaliana* (Halligan and Keightley, 2006; Rajic *et al.*, 2005). The  
39 non-coding genome offers the opportunity to study evolutionary process away from the interference of the direct effects  
40 of natural selection. One such process is the evolution of genomic base composition. The evolution of base content  
41 and its variation within genomes has been the focus of intrigue for many years, such as the question of mammalian  
42 isochore evolution (Eyre-Walker and Hurst, 2001). Genomic GC content is predominately determined by the balance  
43 between the strong (G and C bases) to weak (A and T bases) substitution rate (S→W), in part underpinned by CpG  
44 hypermutability (Hodgkinson and Eyre-Walker, 2011; Hwang and Green, 2004; Ségurel *et al.*, 2014), and the weak to  
45 strong substitution rate (W→S), which is influenced by GC biased gene conversion (gBGC), which favours strong over  
46 weak bases, and is a major determinant of GC content evolution in a broad range of organisms (Bolívar *et al.*, 2016,  
47 2018, 2019; Corcoran *et al.*, 2017; Glémin *et al.*, 2015; Gossmann *et al.*, 2018; Jackson *et al.*, 2017; Muyle *et al.*, 2011;  
48 Ratnakumar *et al.*, 2010; Wallberg *et al.*, 2015). Although, recent experimental based measures of gene conversion in  
49 *Saccharomyces cerevisiae*, *Neurospora crassa*, *Chlamydomonas reinhardtii* and *Arabidopsis thaliana*, did not reveal a  
50 conversion bias (Liu *et al.*, 2018).

51 gBGC is the preferential incorporation of GC alleles over AT alleles during the resolution of heteroduplex DNA  
52 resulting from the repair of double stranded breaks during recombination (Chen *et al.*, 2007; Duret and Galtier, 2009).  
53 This elevates the number of gametes containing GC alleles, as observed in humans (Williams *et al.*, 2015) and birds  
54 (Smeds *et al.*, 2016). As such, gBGC acts to increase the frequency of G and C alleles over A and T alleles, in a manner  
55 that mirrors positive selection (Duret and Galtier, 2009; Galtier and Duret, 2007; Gutz and Leslie, 1976; Nagylaki,  
56 1983). As a result, gBGC is an inconvenient complication when looking for signatures of selection in genomes. For  
57 example, over 20% of identified positively selected genes in the human lineage are possibly just the focus of elevated  
58 gBGC (Ratnakumar *et al.*, 2010). Furthermore, a growing body of literature has demonstrated that gBGC confounds  
59 our ability to estimate parameters such as the rate of adaptation ( $\omega = dN/dS$ ) (Bolívar *et al.*, 2018, 2019; Corcoran  
60 *et al.*, 2017; Gossmann *et al.*, 2018; Ratnakumar *et al.*, 2010; Rousselle *et al.*, 2019) and the proportion of substitutions

61 fixed by positive selection ( $\alpha$ ) (Bolívar *et al.*, 2018; Corcoran *et al.*, 2017; Rousselle *et al.*, 2019). Equally, studying  
62 gBGC in coding regions is inconvenienced by the action of natural selection also acting on those regions, forcing  
63 studies to use putatively neutral sites like third codon positions (Rousselle *et al.*, 2019; Weber *et al.*, 2014) and 4-fold  
64 degenerate sites (Bolívar *et al.*, 2016; Corcoran *et al.*, 2017; Gossmann *et al.*, 2018) reducing the amount of data  
65 available as well as potentially being confounded by codon usage bias (Chamary and Hurst, 2005; Galtier *et al.*, 2018;  
66 Hayes *et al.*, 2020; Jackson *et al.*, 2017; Kunstner *et al.*, 2011).

67 As gBGC is a recombination mediated process, it should co-vary in strength with crossover rate, at different  
68 genomic scales and between species. This is seen in a large body of literature, demonstrating correlations between  
69 recombination rate and GC content (Bolívar *et al.*, 2016; Glémin *et al.*, 2015; Rousselle *et al.*, 2019; Wallberg *et al.*,  
70 2015; Weber *et al.*, 2014), recombination rate and equilibrium GC content (GC\*) (Duret and Arndt, 2008; Muyle  
71 *et al.*, 2011; Singhal *et al.*, 2015), and recombination rate and the population scaled strength of gBGC,  $B = 4N_e b$ ,  
72 where  $N_e$  is the effective population size and  $b$  is the raw strength of conversion bias (Glémin *et al.*, 2015; Wallberg  
73 *et al.*, 2015). However, notably, in *Drosophila* gene conversion rate does not positively correlate with crossover rate  
74 (Comeron *et al.*, 2012). With recombination varying greatly between organisms (Stapley *et al.*, 2017), gBGC can also  
75 be expected to have similar variation in strength and impact. For example, in mammals the recombination landscape  
76 is largely determined by the location of recombination hotspots, determined by the PRDM9 gene (Baudat *et al.*,  
77 2010; Parvanov *et al.*, 2010). This results in areas of greatly elevated recombination rate, and thus strength of gene  
78 conversion relative to background levels, for example, in humans mean  $B$  is estimated at  $\sim 0.4$  (Glémin *et al.*, 2015),  
79 while inside recombination hotspots it reaches as high as  $\sim 18$  (Glémin *et al.*, 2015). In birds, the combination of a  
80 karyotype consisting of a few long macro-chromosomes and many smaller micro-chromosomes (Hansson *et al.*, 2010;  
81 Stapley *et al.*, 2008; van Oers *et al.*, 2014; Zhang *et al.*, 2014) and obligate crossing over causes large chromosomal  
82 differences in recombination rate (Backström *et al.*, 2010; Stapley *et al.*, 2008; van Oers *et al.*, 2014). Additionally,  
83 it has been suggested that birds' lack of PRDM9, has resulted in stable recombination hotspots and conserved  
84 recombination characteristics between species (Singhal *et al.*, 2015). Together this is suggested to allow strong gBGC  
85 to act on the same region of the genome over a longer time period than in mammals (Rousselle *et al.*, 2019; Singhal  
86 *et al.*, 2015), driving GC content increases, with studies reporting that GC content is below GC\* content in most  
87 avian lineages (Bolívar *et al.*, 2016; Mugal *et al.*, 2013; Rousselle *et al.*, 2019; Weber *et al.*, 2014). Furthermore,  
88 some organisms, such as the honey bee *Apis mellifera*, lack pronounced recombination hotspots, yet have very high  
89 genome-wide recombination rate with 5 crossovers per arm and correspondingly elevated mean  $B$  estimates of  $\sim 5$   
90 (Wallberg *et al.*, 2015). Overall, gBGC is seemingly an ubiquitous force with mean  $B$  estimates varying from 0.4 to  
91 5 across the tree of life (Long *et al.*, 2018).

92 As  $B$  is defined as  $4N_e b$ , not only is its strength modulated by recombination rate increasing  $b$  (the strength of  
93 conversion) as outlined above but also by the effective population size ( $N_e$ ). As such species with larger  $N_e$  should  
94 have larger  $B$  and a reduced confounding impact of genetic drift. This has been reported in a few studies, with  
95 correlations between  $N_e$  and GC content at 3rd codon positions (GC3) in birds, largely driven by increased GC in  
96 smaller bodied, larger  $N_e$  species, as well as correlations between  $N_e$  and GC\* (Weber *et al.*, 2014). More recently

97  $B$  at fourfold degenerate sites (4-fold sites) has been shown to correlate with  $N_e$  in great apes (Borges *et al.*, 2019).  
98 However, an analysis of  $B$  more broadly across animal taxa, failed to yield a relationship with  $N_e$  (Galtier *et al.*,  
99 2018). Furthermore, to date the role of  $N_e$  is a less well empirically studied aspect of gBGC and little work has looked  
100 at fine scale variation in the strength of gBGC between species of differing  $N_e$ .

101 The avian system has been the model of choice for many studies addressing GC evolution and biased gene  
102 conversion (Bolívar *et al.*, 2016, 2018, 2019; Corcoran *et al.*, 2017; Gossmann *et al.*, 2018; Rousselle *et al.*, 2019;  
103 Weber *et al.*, 2014). The suitability of avian genomes for addressing these topics stems from their variable intra  
104 genomic recombination landscapes (Backström *et al.*, 2010; Stapley *et al.*, 2008; van Oers *et al.*, 2014) and conserved  
105 recombination hotspots (Singhal *et al.*, 2015) providing a natural experiment for addressing the role of recombination  
106 and  $N_e$  in gBGC and GC content evolution. In addition, birds' conserved karyotype and synteny (Hansson *et al.*,  
107 2010; Stapley *et al.*, 2008; van Oers *et al.*, 2014; Zhang *et al.*, 2014) aids between species comparisons.

108 Of the work on gBGC to date, much has focused on exploring its impact and interaction within genes and  
109 coding regions, largely addressing how it confounds signatures of selection (Bolívar *et al.*, 2019; Corcoran *et al.*,  
110 2017; Gossmann *et al.*, 2018; Ratnakumar *et al.*, 2010; Rousselle *et al.*, 2019). Of those studies that have considered  
111 the action of gene conversion in the non-coding genome (Duret and Arndt, 2008; Glémin *et al.*, 2015; Haddrill and  
112 Charlesworth, 2008; Jackson *et al.*, 2017; Muyle *et al.*, 2011; Wallberg *et al.*, 2015), little work has investigated fine  
113 scale variation within the genome and how this compares between species. Here we investigate variation in the  
114 strength of gBGC within the non-coding genomes of two passerines, the great tit (*Parus major*) and the zebra finch  
115 (*Taeniopygia guttata*), using previously published whole genome resequencing data (Corcoran *et al.*, 2017; Singhal  
116 *et al.*, 2015). We seek to address how conserved the gBGC landscape is between these species and how the strength  
117 of gBGC has been modulated by the recombination rate and  $N_e$  within and between the species.

## 118 **Materials and methods**

### 119 **The dataset**

120 The dataset consisted of 10 European great tits from across the sampling locations in Laine *et al.* (2016), sequenced to  
121 a mean coverage of 44X in Corcoran *et al.* (2017) and 10 zebra finches sequenced to a mean coverage of 22X, a subset  
122 of individuals from the Fowlers Gap population in Australia from the dataset published in Singhal *et al.* (2015). The  
123 dataset is as described in Corcoran *et al.* (2017), but for clarity we will reiterate the main calling pipeline here.

124 SNP calling was performed using GATK v3.4 (Van der Auwera *et al.*, 2013). Raw genotypes were initially called  
125 using the GenotypeGVCF and HaplotypeCaller tools and hard filtered according to the GATK best practice (Van der  
126 Auwera *et al.*, 2013). This call set was used as a training set to perform base quality score recalibration (BQSR).  
127 Variants were then recalled from the recalibrated BAM files both with GATK as above and also using Freebayes v1.02  
128 (Garrison and Marth, 2012). The intersection of the programs' calls was taken and SNPs with less than half, or more  
129 than double the mean depth, and SNPs with a QUAL score less than 20 were removed. This filtered intersection of

130 SNPs was used as a training set to perform variant quality score recalibration (VQSR) on the GATK called variants.  
131 Tranche level thresholds were set at 99% for the zebra finch and 99.9% for the great tit. For both species we obtained  
132 VCF files for SNPs and monomorphic sites from Corcoran *et al.* (2017).

133 Additionally, a three species whole genome alignment between zebra finch (v3.2.4; Warren *et al.*, 2010), great tit  
134 (v1.0.4; Laine *et al.*, 2016) and collared flycatcher (*Ficedula albicollis*) (v1.5; Ellegren *et al.*, 2012) was obtained from  
135 Barton and Zeng (2019), and a three species alignment between chicken (*Gallus gallus*) (v5.0; Hillier *et al.*, 2004),  
136 zebra finch and great tit from Corcoran *et al.* (2017). The former alignment was used to infer the ancestral states of  
137 SNPs, and the latter, with the more distant chicken out-group was used to infer substitution rates and ancestral base  
138 composition (described later). Both of these alignments were generated as follows. Firstly pairwise alignments were  
139 generated with LASTZ (Harris, 2007) between each species and the zebra finch genome, which was used as reference.  
140 These alignments were then chained and netted with axtChain and chainNet respectively (Kent *et al.*, 2003). Single  
141 coverage was ensured for the zebra finch reference genome using single\_cov2.v11 from the MULTIZ package, and  
142 multiple alignments were created from the pairwise alignments using MULTIZ (Blanchette *et al.*, 2004).

## 143 Annotation and filtering

144 We assigned the ancestral states for the SNPs using the whole genome alignment (with collared flycatcher) and  
145 parsimony based approach, where for each species either the reference allele or the alternate allele had to be supported  
146 by both out-groups to be assigned as ancestral.

147 We downloaded the great tit genome annotation (version 1.03) from [ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/001/522/545/GCF\\_001522545.1\\_Parus\\_major1.0.3/GCF\\_001522545.1\\_Parus\\_major1.0.3\\_genomic.gff.gz](ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/001/522/545/GCF_001522545.1_Parus_major1.0.3/GCF_001522545.1_Parus_major1.0.3_genomic.gff.gz) (last  
148 accessed 05/03/19) and the zebra finch annotation from [ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/151/805/GCF\\_000151805.1-Taeniopygia\\_guttata-3.2.4](ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/151/805/GCF_000151805.1-Taeniopygia_guttata-3.2.4) (last accessed on 05/03/19). We used the annotations to re-  
149 move variants falling within exons. Additionally coordinates for ultra-conserved non-coding elements (UCNEs) in the  
150 zebra finch genome (taeGut1) were obtained from [ftp://cgg.vital-it.ch/UCNEbase/custom\\_tracks\\_UCSC/UCNEs\\_taeGut1.bed](ftp://cgg.vital-it.ch/UCNEbase/custom_tracks_UCSC/UCNEs_taeGut1.bed) (last accessed 05/03/19). We identified the corresponding positions in the great tit in the whole genome  
151 alignment, before removing any variants falling within UCNEs. Additionally we restricted our analysis to the auto-  
152 somes, removing the Z chromosome. This left non-coding datasets of putatively neutral variants, numbering 9,800,315  
153 SNPs for great tit, and 29,973,954 SNPs for zebra finch.

154 From our non-coding SNP dataset we generated an additional subset, with CpG sites excluded, where a CpG site  
155 was defined as any site where at least one of the alleles of the site was in a 5' → 3' CpG dinucleotide or in a 3' → 5'  
156 GpC dinucleotide.

## 160 Orthologous window preparation

161 The zebra finch genome was divided into 1Mb non-overlapping windows and we used the three species whole genome  
162 alignment (zebra finch, great tit, collared flycatcher) to identify the aligned sequence and coordinates in the great tit

163 genome and extracted variants and numbers of callable sites from our VCF files. For each window in each species  
164 we calculated the GC content using the respective reference genomes. GC content was calculated for all sites in the  
165 window, and for non-CpG sites. Secondly, we calculated crossover rate for each window, using the available linkage  
166 map data for each species (Stapley *et al.*, 2008; van Oers *et al.*, 2014, for zebra finch and great tit respectively) and  
167 the pipeline outlined in Corcoran *et al.* (2017).

## 168 Estimating the strength of gene conversion

169 We extracted the number of callable sites for weak bases (A and T nucleotides) and strong bases (G and C nucleotides)  
170 along with the site frequency spectra for weak to strong mutations (*WS*), strong to weak mutations (*SW*) and weak  
171 to weak and strong to strong mutations (*WWSS*) in all windows and datasets. We then applied the *M1\** model of  
172 Glémin *et al.* (2015), implemented in the **anavar** package (Barton and Zeng, 2018), to all windows with at least 1,000  
173 SNPs. Briefly, the model estimates the population scaled mutation rate ( $\theta = 4N_e\mu$ ), the population scaled strength of  
174 gBGC ( $B = 4N_e b$ ) and estimates and controls for polarisation error for both *SW* and *WS* mutations using *WWSS*  
175 sites as a neutral reference unaffected by gBGC. Demography is controlled for using the method of Eyre-Walker *et al.*  
176 (2006), which has been shown previously to obtain similar results to a method that explicitly model recent changes  
177 in population size (Jackson *et al.*, 2017).

178 We performed multiple regressions in R (R Core Team, 2015) to estimate the relative contributions of crossover  
179 rate and local  $N_e$  (using nucleotide diversity  $[\pi]$  as a measure of  $N_e$ ) in determining  $B$ , we ran these analysis using  
180 crossover rate and separately, GC content as measures of recombination rate. We estimated the relative importance of  
181 the predictors (as a proportion of the total variance explained) using the ‘pmvd’ method implemented in the **relaimpo**  
182 package (Groemping, 2006).

## 183 Equilibrium GC content

184 We estimated the ancestral GC content per window for the lineage leading to great tits and zebra finches using the  
185 whole genome alignment (containing chicken, zebra finch and great tit) and the GTR-NH<sub>b</sub> model in baseml within  
186 PAML (Yang, 2007). The model allows for non-stationary base content and for independent substitution rates on each  
187 branch. From the model we obtained the posterior probabilities of the ancestral states and weighted each ancestral  
188 nucleotide by this probability (as in Matsumoto *et al.*, 2015) to reconstruct ancestral GC content with uncertainty  
189 incorporated. We then estimated the rate of *WS* substitutions

$$r_{WS} = \frac{n_{WS}}{n_W} \quad (1)$$

190 where  $n_{WS}$  is the number of *WS* substitutions and  $n_W$  is the number of weak bases (As and Ts) in the ancestral  
191 sequence. Similarly we estimated the rate of *SW* substitutions

$$r_{SW} = \frac{n_{SW}}{n_S} \quad (2)$$

192 where  $n_{SW}$  is the number of  $SW$  substitutions and  $n_S$  is the number of strong (Gs and Cs) bases in the ancestral  
193 sequence. Finally we estimated the equilibrium GC content

$$GC^* = \frac{r_{WS}}{r_{WS} + r_{SW}} \quad (3)$$

194 The GTR-NH<sub>6</sub> model was a better fit than the GTR model, which assumes base composition is at equilibrium, for  
195 all but five windows as judged by likelihood ratio tests (data not shown). Additionally, the model estimates of  $GC_{div}^*$   
196 correlated strongly with those derived from parsimony estimates of the substitution rates for both great tit (Pearson's  
197  $r = 0.94$ ,  $p < 2.2 \times 10^{-16}$ ) and zebra finch (Pearson's  $r = 0.96$ ,  $p < 2.2 \times 10^{-16}$ ), although the mean  $GC_{div}^*$  was lower  
198 for the model estimates than the parsimony estimates in both species (0.39 versus 0.43 respectively for great tit and  
199 0.38 versus 0.42 respectively for zebra finch).

200 To obtain a more recent view of the base composition evolution and gBGC we also calculated  $GC_{pol}^*$  from our  
201 application of the Glémin *et al.* (2015) model to our polymorphism dataset. In order to do so we took the estimates  
202 of  $B$  ( $B = 4N_e b$ ) and mutation rates ( $\theta = 4N_e \mu$ ) estimated per window by `anavar` and substituted them into

$$r_{ij} = \theta_{ij} \frac{B_{ij}}{1 - e^{-B_{ij}}} \quad (4)$$

203 where  $r_{ij}$  is the fixation rate of mutations from  $i$  to  $j$  and where  $B_{WS} = -B_{SW} = B$ . The resulting fixation rates  
204 were then substituted into equation 3 to obtain  $GC^*$ .

## 205 Demographic analysis

206 To investigate the demographic history in the zebra finch and the great tit we fitted demographic models to the data  
207 using the `VarNe` package Zeng *et al.* (2019). The package performs maximum likelihood estimation of a number of  
208 population genetic parameters, including  $\theta$  ( $4N_e \mu$ ), the magnitude of a population size change ( $g$ ), the timing of  
209 the event ( $\tau$ , in units of  $2N_e$ ) and the rate of ancestral state misidentification ( $\epsilon$ ), allowing population size changes  
210 between a specified number of time points, or epochs, from the site frequency spectrum of a target locus. We applied  
211 1 epoch and 2 epoch models to the summed site frequency spectra for  $WWSS$  (GC conservative) non-coding SNPs  
212 from our window dataset. We tested whether the 2 epoch model (variable population size) was a better fit than the  
213 1 epoch model (constant size), using likelihood ratio tests. We performed 100 rounds of bootstrapping by resampling  
214 windows from our window dataset with replacement.

215 We also applied the 2 epoch model above individually to each window in our dataset to obtain local estimates of  
216 the magnitude of  $N_e$  change. For these analyses we required windows to have a minimum of 1000 SNPs and windows  
217 that failed to return reliable parameter estimates were excluded (67 windows in the great tit, 4 windows in the zebra  
218 finch).

219 In order to infer how much our polymorphism based estimate of the equilibrium GC content ( $GC_{pol}^*$ ) might differ  
220 prior to the inferred population size change in each species, we divided our estimates of  $B_{WS}$ ,  $\theta_{WS}$ ,  $B_{SW}$  and  $\theta_{SW}$   
221 by a correction factor  $C$ , as a function of  $g$  and  $\tau$  estimates per window:

$$C = g + (1 - g)e^{-\tau/g} \quad (5)$$

222 We then substituted the rescaled values into equation 4, to calculate the fixation probabilities for *WS* and *SW*  
223 polymorphisms under the reduced *B* scenario. The fixation probabilities were then substituted into equation 3 to  
224 calculate  $GC^*$ .

## 225 Data availability

226 All scripts and command lines used in the analysis pipeline can be found at: [https://github.com/henryjuho/biased\\_](https://github.com/henryjuho/biased_gene_conversion)  
227 [gene\\_conversion](#). The VCF files, whole genome alignments and orthologous window coordinates are accessible at:  
228 [link](#).

## 229 Results

### 230 Summary of the window dataset

231 We used a whole genome alignment between zebra finch, great tit and collared fly catcher (*Ficedula albicollis*) to  
232 identify 1Mb orthologous windows between the zebra finch and great tit. This resulted in 904 1Mb windows in zebra  
233 finch genome and 898 orthologous windows in the great tit genome (table 1). The lower number of great tit windows  
234 is due to gaps in the whole genome alignment. We used the respective genome annotations to identify non-coding  
235 regions within these windows, in which we identified single nucleotide polymorphisms (SNPs) using a resequencing  
236 dataset of 10 zebra finches (from Singhal *et al.*, 2015) and 10 great tits (from Corcoran *et al.*, 2017). This resulted  
237 in similar numbers of callable sites in both species, roughly 500,000 bp per 1 Mb window; this drop is a result of our  
238 focus on non-coding regions (excluding ultra-conserved non-coding elements [UCNEs]), and our maximum parsimony  
239 approach to assigning ancestral states, which is dependant on coverage of all species in our whole genome alignment  
240 and no ambiguity between out-groups. When considering variants per window, we see that the mean number of  
241 variants is higher in the zebra finch, consistent with a larger effective population size in the zebra finch (Corcoran  
242 *et al.*, 2017). We see very similar mean GC content and mean crossover rates in both species, with strong correlations  
243 between the two species' GC content (Pearson's  $r = 0.83$ ,  $p = 1.6 \times 10^{-230}$ , figure S1a) and crossover rate (Spearman's  
244  $\rho = 0.72$ ,  $p = 2.6 \times 10^{-140}$ , figure S1b) across the dataset, as well as positive correlations between GC content and  
245 crossover rate within each species (great tit: Spearman's  $\rho = 0.57$ ,  $p = 3.8 \times 10^{-79}$ , zebra finch: Spearman's  $\rho = 0.53$ ,  
246  $p = 4.2 \times 10^{-67}$ , figure S2).

### 247 The strength of gene conversion correlates with crossover rate and $N_e$

248 To estimate the population scaled strength of gBGC (*B*), we applied the Glémin *et al.* (2015) model to each window  
249 in our dataset. The resulting estimates of *B* positively correlate with both crossover rate and  $\pi$  (as a proxy for local



Table 1: Summary of the window dataset, showing means and the 2.5 and 97.5 percentiles in brackets. Crossover rates are log10 transformed.

Measure	great tit	zebra finch
windows	898	904
callable sites	523858 (21580, 726488)	498785 (79743, 711346)
$n_{SNP}$	5895 (239, 9766)	21321 (989, 37847)
GC content	0.41 (0.34, 0.51)	0.41 (0.35, 0.51)
Crossover rate (cM/Mb)	0.48 (0, 0.97)	0.41 (0, 0.96)

Table 2: Results of multiple regression analysis of the strength of gene conversion ( $B$ ) against GC content and  $\pi$ , and against crossover rate and  $\pi$ , separately, for both species. Importance is the relative importance (as a proportion of the total variance explained) as estimated using the pmvd method implemented in the `relaimpo` package (Groemping, 2006).

model	species	variable	estimate	importance	p value	$R^2$
$B \sim \log_{10}(\text{crossover rate} + 1) + \pi$	great tit	crossover rate	0.652	0.94	$< 2 \times 10^{-16}$	0.264
		$\pi$	59.2	0.06	$5.86 \times 10^{-5}$	
	zebra finch	crossover rate	1.46	0.81	$< 2 \times 10^{-16}$	0.505
		$\pi$	78.6	0.19	$< 2 \times 10^{-16}$	
$B \sim \text{GC content} + \pi$	great tit	GC content	4.90	0.88	$< 2 \times 10^{-16}$	0.371
		$\pi$	105	0.12	$9.73 \times 10^{-15}$	
	zebra finch	GC content	13.0	0.94	$< 2 \times 10^{-16}$	0.656
		$\pi$	50.0	0.06	$3.86 \times 10^{-16}$	

250  $N_e$ , allowing us to separate the contributions of  $N_e$  to the compound parameter  $B = 4N_e b$  in both the great tit  
 251 and the zebra finch (table 2, figure 1). The relationships are stronger when using mean GC content as a proxy for  
 252 recombination rate in both species (table 2, figure S3) and all relationships are maintained when performed on a  
 253 dataset filtered for CpG sites (table S1). Crossover rate or mean GC content explains a larger proportion of the total  
 254 variance (80 – 95%) than  $\pi$  within both species (table 2, table S1).

## 255 $B$ is correlated between the species

256 Comparison of the model estimates of  $B$  between zebra finch and great tit show a significantly larger mean  $B$  value  
 257 in zebra finch ( $\bar{B} = 0.90$ ) than great tit ( $\bar{B} = 0.40$ ) (Wilcoxon rank sum,  $W = 491903$ ,  $p = 2.5 \times 10^{-49}$ ; figure  
 258 2a), inline with the species' twofold difference in  $N_e$  (Corcoran *et al.*, 2017). However, when we standardise our  $B$   
 259 estimates by  $\pi$  as a measure of  $N_e$ , the difference between the two species is greatly reduced and the distributions  
 260 of  $B/\pi$  are similar in both species (figure 2b). However,  $B/\pi$  is slightly, but significantly larger in the great tit  
 261 ( $\bar{B}/\pi = 118.2$  and  $80.8$  for great tit and zebra finch respectively, Wilcoxon rank sum,  $W = 305880$ ,  $p = 6.1 \times 10^{-10}$ ).  
 262 We also see a positive correlation between the ratio of the species' nucleotide diversity ( $\pi_{zf}/\pi_{gt}$ ) and the ratio of  
 263 the species'  $B$  ( $B_{zf}/B_{gt}$ ) (Spearman's  $\rho = 0.44$ ,  $p < 2.2 \times 10^{-16}$ ), supporting the idea that  $N_e$  drives the between  
 264 species differences in  $B$ . Furthermore, we see a strong correlation between  $B$  in the great tit and  $B$  in the zebra finch  
 265 (Pearson's  $r = 0.50$ ,  $p < 2.2 \times 10^{-16}$ , figure 3) as well as between  $B/\pi$  in great tit and  $B/\pi$  in zebra finch (Pearson's  
 266  $r = 0.38$ ,  $p < 2.2 \times 10^{-16}$ ), in keeping with the conserved crossover rate and GC content between species reported  
 267 above.

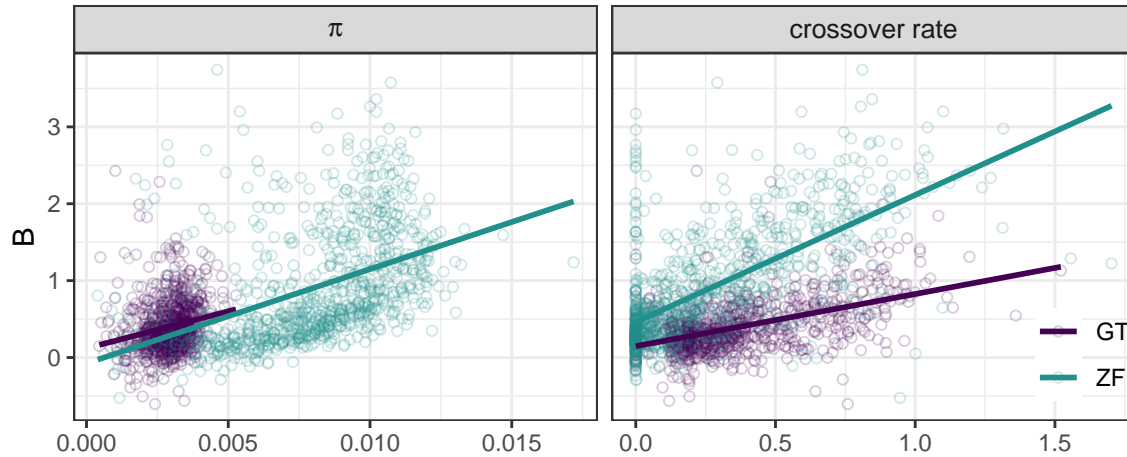


Figure 1: The relationship between nucleotide diversity ( $\pi$ ) and the strength of gene conversion ( $B$ ) (left panel) and mean window crossover rate and  $B$  (right panel) in the great tit (purple) and zebra finch (turquoise). Multiple regression results can be seen in table 2.

## Equilibrium GC content

To assess the longer term GC dynamics of both the great tit and zebra finch genomes, we calculated the equilibrium GC content ( $GC^*$ ), which is the GC content that when reached will result in equal numbers of GC alleles fixed as lost.

Firstly, we calculated  $GC^*$  using divergence data ( $GC_{div}^*$ ) for each lineage, using the WS and SW substitution rates estimated in PAML (see methods). This provides a long term average of  $GC^*$  since the two species diverged. This gave a mean  $GC_{div}^*$  of 0.39 for great tit and 0.38 for zebra finch, both of which are similar to, but significantly below, the mean GC contents in our alignment datasets of 0.40 for both great tit (Wilcoxon rank sum,  $W = 282790$ ,  $p = 1.1 \times 10^{-8}$ ) and zebra finch (Wilcoxon rank sum,  $W = 241190$ ,  $p < 2.2 \times 10^{-16}$ ) (figure 4). Note the alignment dataset is a subset of the main dataset (as coverage is required across all species in the chicken/zebra finch/great tit alignment) and yields slightly lower mean GC than reported in table 1.  $B$  positively correlates with  $GC_{div}^*$  in both great tit (Pearson's  $r = 0.54$ ,  $p < 2.28 \times 10^{-55}$ ) and zebra finch (Pearson's  $r = 0.81$ ,  $p < 8.22 \times 10^{-181}$ ). Similar relationships are seen between  $GC_{div}^*$  and crossover rate (Spearman's  $\rho = 0.55$ ,  $p = 6.02 \times 10^{-62}$  for great tit and Spearman's  $\rho = 0.66$ ,  $p = 3.85 \times 10^{-98}$  for zebra finch) and between  $GC_{div}^*$  and current GC content (Pearson's  $r = 0.56$ ,  $p = 9.32 \times 10^{-65}$  for great tit and Pearson's  $r = 0.77$ ,  $p = 1.49 \times 10^{-148}$  for zebra finch).

Secondly, to look at base composition evolution over a more recent time scale we also calculated  $GC^*$  from polymorphism data, using our  $\theta$  and  $B$  estimates derived from the Glémin *et al.* (2015) model (see methods), henceforth  $GC_{pol}^*$ . This approach yielded markedly higher equilibrium GC content estimates than the substitution rate based approach, for both great tit (Wilcoxon rank sum,  $W = 518421$ ,  $p = 1.48 \times 10^{-225}$ ,  $\bar{GC}_{pol}^* = 0.63$ ) and zebra finch (Wilcoxon rank sum,  $W = 575196$ ,  $p = 1.24 \times 10^{-245}$ ,  $\bar{GC}_{pol}^* = 0.72$ ).

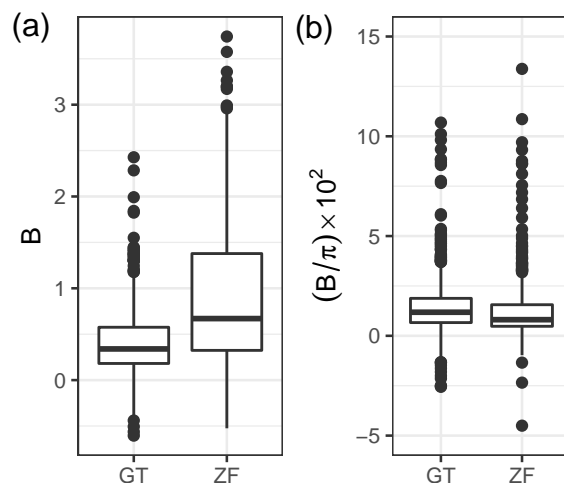


Figure 2: Comparison of the distribution of  $B$  values (population scaled strength of biased gene conversion) (a) and  $B$  standardised by  $\pi$  as a proxy for the effective population size  $N_e$  (b) between the great tit (GT) and zebra finch (ZF). The y axis for b has been cropped for clarity.

## Evidence of population expansions

In order to understand the effects of recent demographic changes on the difference between our longer term measures of  $GC_{div}^*$  and our more recent  $GC_{pol}^*$  estimates, we fitted demographic models to each species using the **VarNe** package (Zeng *et al.*, 2019). The models estimate the magnitude ( $g$ ) and timing of population size changes ( $\tau$ , in units of  $2N_e$ ) between different time points or ‘epochs’. In both the zebra finch and the great tit a 2 epoch model (table S3) fit the data significantly better than a 1 epoch model (i.e. a model with constant population size) as judged by likelihood ratio tests. For the zebra finch we estimate a  $g$  of 12.3 and  $\tau$  of 1.25, suggesting a large population expansion  $\sim 495$  thousand years ago (table S3). In the great tit we see lower values with a  $g$  of 1.89 and  $\tau$  of 0.208, characterising a smaller, more recent population expansion  $\sim 140$  thousand years ago (table S3).

## Local $N_e$ increase correlates with increases in equilibrium GC content in the zebra finch

Nucleotide diversity is positively correlated with recombination rate in both the great tit and zebra finch (Corcoran *et al.*, 2017), showing  $N_e$  varies locally within their genomes. As loci with differing  $N_e$  can respond differently to a shared demographic change (see Zeng *et al.*, 2019), we sort to investigate how historical changes in local  $N_e$  have impacted equilibrium GC content, and the difference between our  $GC_{div}^*$  and  $GC_{pol}^*$  estimates. In each species, we refitted the ‘2 epoch’ model in **VarNe**, to each window in our orthologous window dataset. The mean maximum likelihood parameter estimates across all windows agreed with those from the model fitted to the dataset as a whole, although were slightly higher, probably a result of our requirement of a minimum of 1000 SNPs per window to provide

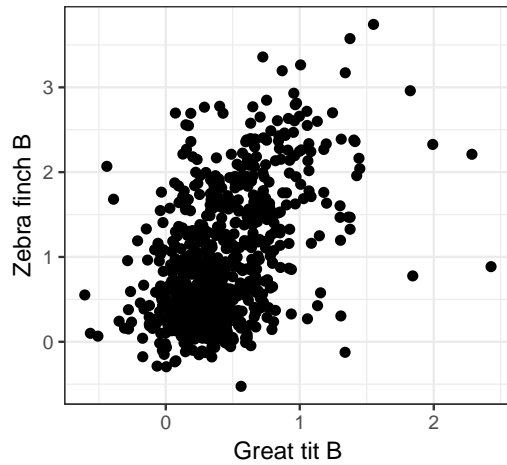


Figure 3: The strength of biased gene conversion ( $B$ ) in the zebra finch positively correlates with  $B$  in the great tit.

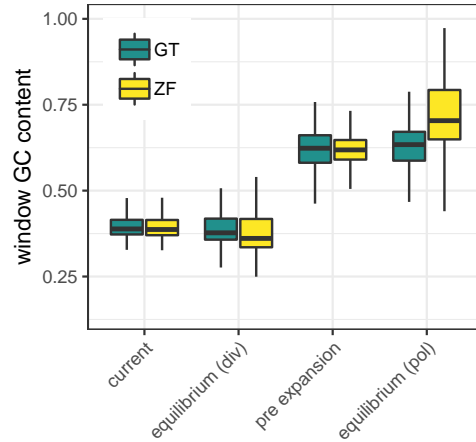


Figure 4: Per window estimates of current GC content, equilibrium GC content from both divergence data (div) and polymorphism data (pol) and estimates of equilibrium GC content before expansion (pre expansion) for both species.

306 sufficient power, excluding the lowest  $N_e$  windows (table S4).

307 For each window, we divided our estimates of  $\theta$  and  $B$  from the Glémin *et al.* (2015) model by a rescaling factor  $C$   
 308 ( $\bar{C}_{gt} = 1.18$  and  $\bar{C}_{zf} = 2.32$ ), a function of the window's  $g$  and  $\tau$  estimates (see equation 5) to control for the effects  
 309 of recent population expansion. We used these rescaled values to obtain per window estimates of  $GC^*$ , prior to the  
 310 inferred local  $N_e$  increases. This approach yielded a mean pre-expansion  $GC^*$  of 0.62 in both the zebra finch and the  
 311 great tit (figure 4), demonstrating the transient effect of recent population size changes on equilibrium GC content.  
 312 These  $GC^*$  estimates are still high relative to  $GC^*_{div}$ , potentially due to population expansions taking place before  
 313 the most recent common ancestor of the polymorphism samples.

314 Additionally, we compared the per window values of  $C$  (a measure of the impact of  $N_e$  increase on  $B$ ) with the  
 315 difference between our two  $GC^*$  estimates. This returned a significant positive correlation between  $GC^*$  increase  
 316 ( $GC^*_{pol} - GC^*_{div}$ ) and  $C$  in the zebra finch (Spearman's  $\rho = 0.46$ ,  $p < 2.2 \times 10^{-16}$ ) and a weak positive correlation  
 317 in the great tit (Spearman's  $\rho = 0.081$ ,  $p = 0.036$ ). The stronger correlation in the zebra finch is consistent with an  
 318 older and larger expansion in this species providing more time for evolution to influence  $C$  and  $GC^*$ .

## Discussion

Most contemporary studies on the role of GC biased gene conversion (gBGC) in genome evolution have focused on coding regions where gBGC is confounded by selection, (Bolívar *et al.*, 2019; Corcoran *et al.*, 2017; Gossmann *et al.*, 2018; Ratnakumar *et al.*, 2010; Rousselle *et al.*, 2019) and processes like codon usage bias (Haddrill *et al.*, 2008; Jackson *et al.*, 2017). Additionally few of these studies have looked at the impact of  $N_e$  on the strength of gBGC. Here we analyse re-sequencing data for 10 great tits (Corcoran *et al.*, 2017), and 10 zebra finches (Singhal *et al.*, 2015). Using non-overlapping 1Mb orthologous windows, we investigate how the strength and impact of gBGC varies both within and between the non-coding genomes of these birds.

### The strength of gene conversion is modulated by $N_e$

Our mean estimates of  $B$  in the great tit and the zebra finch of 0.40 and 0.90 respectively, are similar to mean genome wide estimates of  $B$  in humans of 0.38 (Glémin *et al.*, 2015), and fall at the lower end of the  $B$  range of 0.4 to 5 reported by Long *et al.* (2018) in a comparative study with taxa from across the tree of life. Mutations with  $N_e s < 1$  (here  $B = 4N_e b < 4$ ) are considered effectively neutral, our mean  $B$  estimates fall below 1, suggesting gBGC in the non-coding regions of these species is operating at low efficiency.

$B$  correlates with both recombination rate and  $\pi$  in these species (table 2) suggesting both parameters are modulating  $B$  in their genomes, although recombination rate has the larger impact (when measured by crossover rate or mean GC content), particularly in the great tit. This is consistent with elevated gBGC in regions with higher recombination rate in humans (Glémin *et al.*, 2015) and correlations between GC content at 4-fold sites and recombination rate in flycatchers (Bolívar *et al.*, 2016), although these analyses did not control for local  $N_e$ . When using GC content as a measure of recombination rate instead of crossover rate these relationships are strengthened. This may reflect that GC content is a better measure of long term recombination rate, that our crossover rate estimates are constrained by the density of the linkage maps available (Stapley *et al.*, 2008; van Oers *et al.*, 2014), lower variance in our GC estimates (table 1), or a mixture of the three.

The conservation of the biased gene conversion landscape between the zebra finch and great tit, as seen by the strong correlation of window  $B$  estimates between the species, is relatively intuitive with GC content and crossover rate also correlating well between the species and likely a result of birds' conserved recombination hotspots (Singhal *et al.*, 2015), karyotype and synteny (Hansson *et al.*, 2010; Stapley *et al.*, 2008; van Oers *et al.*, 2014; Zhang *et al.*, 2014). Consistently, we also see similar mean crossover rates in each species (table 1).

Nonetheless, mean  $B$  is approximately twofold higher in the zebra finch. As  $B$  is the product of  $b$  (the strength of biased gene conversion) and  $N_e$ , either parameter could be driving this increase. When we standardise  $B$  by  $\pi$  (as a measure of  $N_e$ ), the between species difference is greatly reduced. This, combined with the correlation of the ratios of between species  $B$  ( $B_{zf}/B_{gt}$ ) and  $\pi$  ( $\pi_{zf}/\pi_{gt}$ ), suggests the twofold larger  $N_e$  in the zebra finch (Corcoran *et al.*, 2017) is elevating its  $B$ . This also implies that  $b$  is comparable between the species and has remained relatively stable since their divergence. Consistently, GC3 content correlates with  $N_e$  (using life history traits as proxies) across

353 the avian phylogeny (Weber *et al.*, 2014). More broadly, it fits with findings in great apes, where  $B$  at 4-fold sites  
354 correlates with  $N_e$  (Borges *et al.*, 2019) and amongst rice species (*Oryza spp.*), where selfing species (with reduced  
355  $N_e$ ) also have lower  $B$  estimates (Muyle *et al.*, 2011). However, a recent analysis by Galtier *et al.* (2018) between  
356 more diverged species failed to find a relationship between  $B$  and  $N_e$ , with the authors suggesting that  $b$  may be  
357 inversely related to  $N_e$  between distant taxa, and only remain homogenous within groups, such as birds, suggesting  
358 that  $B$  only responds to  $N_e$  over small time-scales.

## 359 Non-coding equilibrium GC content

360 Our two measures of equilibrium GC content ( $GC^*$ , the theoretical GC content at which the same number of GC  
361 alleles are fixed as AT alleles, and thus stable GC content is reached), from divergence data ( $GC_{div}^*$ ) and polymorphism  
362 data ( $GC_{pol}^*$ ), respectively provide a longer term and more recent insight into  $GC^*$ .

363  $GC_{div}^*$  is similar to, albeit significantly lower than, current GC content in both species. This is at odds with previous  
364 avian studies where GC content is below  $GC^*$  in most lineages (Bolívar *et al.*, 2016; Rousselle *et al.*, 2019; Weber  
365 *et al.*, 2014). However, these studies focus on coding regions, which are have elevated GC content and recombination  
366 rates over non-coding regions in birds (Singhal *et al.*, 2015; Weber *et al.*, 2014); in our dataset, GC content is  $\sim 10\%$   
367 higher in coding regions than non-coding regions (table S2). Consequently, gBGC is likely stronger in coding regions,  
368 as suggested by  $GC_{div}^*$  estimates of 0.6 – 0.8 at fourfold sites in collared flycatcher (Bolívar *et al.*, 2016) and a median  
369  $GC_{div}^*$  of 0.6 at 3rd codon positions across 48 bird species (Weber *et al.*, 2014), compared to our non-coding  $GC_{div}^*$  of  
370 0.39 in the great tit and 0.38 in the zebra finch. These differing dynamics may be contributed to by the avian micro-  
371 chromosomes which are characterised by high gene density, and high recombination rates stemming from obligate  
372 crossing over and their short length (Burt, 2002; Stapley *et al.*, 2008; van Oers *et al.*, 2014). Equally, if codon usage  
373 bias (CUB) is operating in addition to gBGC (de Procé *et al.*, 2012; Galtier *et al.*, 2018) and favours G and C ending  
374 codons (de Procé *et al.*, 2012) this could elevate avian coding GC over non-coding GC, and also inflate estimates of  
375 gBGC in coding regions, however, evidence for CUB in birds is lacking. Overall, it seems these regions have been  
376 evolving towards different equilibria, similar to some species of rice (Muyle *et al.*, 2011), with weak gBGC allowing  
377 for more AT biased fixation patterns (see McVean and Charlesworth, 1999) and a slightly decreasing GC content in  
378 non-coding regions since the great tit zebra finch split.

## 379 The effect of demography on $B$ and $GC^*$

380 Our mean  $GC_{pol}^*$  estimates are higher than our  $GC_{div}^*$  values, 0.63 versus 0.39 for great tit and 0.72 versus 0.38 for  
381 zebra finch.  $GC_{div}^*$  represents a long term average of  $GC^*$  since the divergence of the great tit and zebra finch lineages  
382 40 to 45 million years ago (Barker *et al.*, 2004), whereas  $GC_{pol}^*$  provides a more recent snapshot, of the order of  $4N_e$   
383 generations ago, around  $\sim 3.5$  and  $\sim 4.3$  million years ago for the great tit and zebra finch respectively (estimated  
384 using the current  $N_e$  estimates and generation times in table S5). Consequently, our higher  $GC_{pol}^*$  estimates suggests  
385  $B$  is currently higher than the long term average for the species, this is the opposite to what is seen in *Drosophila*

386 *melanogaster*, where longer term estimates of  $B$  are higher than those from the Glémin *et al.* (2015) model (Jackson  
387 *et al.*, 2017). As  $B$  is the product of  $b$  (the underlying strength of conversion bias) and  $N_e$ , this increase could be  
388 driven by increases in the population size and/or  $b$  through changing recombination rates. As recombination rates  
389 are relatively stable and conserved in these species (Singhal *et al.*, 2015; van Oers *et al.*, 2014, this study), it seems  
390 more probable the current elevation of  $B$  is driven by changes in  $N_e$ .

391 Here, we estimate  $\sim 12$ -fold and  $\sim 2$ -fold population expansions for the zebra finch and great tit respectively, in  
392 agreement with previous evidence for expansions in both species (Balakrishnan and Edwards, 2008; Corcoran *et al.*,  
393 2017; Laine *et al.*, 2016). The magnitude of the great tit expansion is similar to reported values of 2.75 (Laine *et al.*,  
394 2016), 2.31 (Corcoran *et al.*, 2017) and 1.68 (Hayes *et al.*, 2020). The zebra finch expansion magnitude of 12.3 is close  
395 to the estimate of 10 from Corcoran *et al.* (2017), the upper limit of the method used. The larger increase in  $N_e$  for  
396 the zebra finch is consistent with the greater difference in  $GC^*$  measures in this species (figure 4). Furthermore, our  
397 estimates of  $GC_{pol}^*$  corrected for the inferred population expansions are 0.62 in both species, suggesting each species'  
398 average  $N_e$  have remained similar since they diverged. The difference between  $GC_{pol}^*$  and  $GC_{div}^*$  is reduced by 29%  
399 after correction in the zebra finch, but only by 4% in the great tit. Concordantly, the difference between  $GC_{pol}^*$  and  
400  $GC_{div}^*$  correlates well with our correction factor  $C$ , a measure of the impact of  $N_e$  increase, in zebra finch only. As the  
401 polymorphism data spans at most 10% of the species divergence time, most of the demographic history since their  
402 split is not captured in our analysis, thus the modest impact of the recent expansions on  $GC^*$  is perhaps unsurprising.

## 403 Conclusion

404 We show that the underlying strength of gene conversion  $b$  is conserved between the great tit and zebra finch, with  
405 the zebra finch's larger population scaled strength of gBGC,  $B$ , due to its larger effective population size. Within  
406 each species' genome, variation in  $B$  is driven by variation in both recombination rate and local  $N_e$ , with the former  
407 having the larger impact.

408 When considering the equilibrium GC content, we see that  $GC_{div}^*$  and  $GC^*$  prior to the inferred population  
409 expansions are similar between the great tit and zebra finch, suggesting that they have had similar average  $N_e$  since  
410 their divergence. Our higher  $GC_{pol}^*$  estimates are likely explained by the short timescale covered by the polymorphism  
411 data relative to the divergence data.

## 412 Acknowledgements

413 We thank Pádraic Corcoran for providing an initial implementation of the window pipeline and Brian Charlesworth for  
414 his comments on the manuscript. This work was supported by a PhD studentship funded by the Department of Animal  
415 and Plant Sciences, University of Sheffield, to H.J.B. Support was also provided by the Natural Environment Research  
416 Council via a research grant awarded to K.Z. (NE/L005328/1). The analyses were performed on the University of  
417 Sheffield's high performance computing cluster 'ShARC'.

## 418 References

- 419 Backström, N., Forstmeier, W., Schielzeth, H., Mellenius, H., Nam, K., Bolund, E., Webster, M. T., Ost, T., Schneider,  
420 M., Kempenaers, B., and Ellegren, H. 2010. The recombination landscape of the zebra finch *Taeniopygia guttata*  
421 genome. *Genome Res*, 20(4): 485–95.
- 422 Balakrishnan, C. N. and Edwards, S. V. 2008. Nucleotide Variation, Linkage Disequilibrium and Founder-Facilitated  
423 Speciation in Wild Populations of the Zebra Finch (*Taeniopygia guttata*). *Genetics*, 181(2): 645–660.
- 424 Barker, F. K., Cibois, A., Schikler, P., Feinstein, J., and Cracraft, J. 2004. Phylogeny and diversification of the largest  
425 avian radiation. *Proceedings of the National Academy of Sciences*, 101(30): 11040–11045.
- 426 Barton, H. J. and Zeng, K. 2018. New Methods for Inferring the Distribution of Fitness Effects for INDELS and  
427 SNPs. *Molecular Biology and Evolution*, 35(6): 1536–1546.
- 428 Barton, H. J. and Zeng, K. 2019. The Impact of Natural Selection on Short Insertion and Deletion Variation in the  
429 Great Tit Genome. *Genome Biology and Evolution*, 11(6): 1514–1524.
- 430 Baudat, F., Buard, J., Grey, C., Fledel-Alon, A., Ober, C., Przeworski, M., Coop, G., and Massy, B. d. 2010. PRDM9  
431 Is a Major Determinant of Meiotic Recombination Hotspots in Humans and Mice. *Science*, 327(5967): 836–840.
- 432 Blanchette, M., Kent, W. J., Riemer, C., Elnitski, L., Smit, A. F. A., Roskin, K. M., Baertsch, R., Rosenbloom, K.,  
433 Clawson, H., Green, E. D., Haussler, D., and Miller, W. 2004. Aligning Multiple Genomic Sequences With the  
434 Threaded Blockset Aligner. *Genome Research*, 14(4): 708–715.
- 435 Bolívar, P., Mugal, C. F., Nater, A., and Ellegren, H. 2016. Recombination rate variation modulates gene sequence  
436 evolution mainly via GC-biased gene conversion, not Hill–Robertson interference, in an avian system. *Molecular*  
437 *biology and evolution*, 33(1): 216–227.
- 438 Bolívar, P., Mugal, C. F., Rossi, M., Nater, A., Wang, M., Dutoit, L., and Ellegren, H. 2018. Biased Inference of  
439 Selection Due to GC-Biased Gene Conversion and the Rate of Protein Evolution in Flycatchers When Accounting  
440 for It. *Molecular Biology and Evolution*, 35(10): 2475–2486.
- 441 Bolívar, P., Guéguen, L., Duret, L., Ellegren, H., and Mugal, C. F. 2019. GC-biased gene conversion conceals the  
442 prediction of the nearly neutral theory in avian genomes. *Genome Biology*, 20(1): 5.
- 443 Borges, R., Szöllösi, G. J., and Kosiol, C. 2019. Quantifying GC-Biased Gene Conversion in Great Ape Genomes  
444 Using Polymorphism-Aware Models. *Genetics*, 212(4): 1321–1336.
- 445 Burt, D. W. 2002. Origin and evolution of avian microchromosomes. *Cytogenetic and Genome Research*, 96(1-4):  
446 97–112.
- 447 Chamary, J. and Hurst, L. D. 2005. Evidence for selection on synonymous mutations affecting stability of mRNA  
448 secondary structure in mammals. *Genome Biology*, 6(9): R75.



- 449 Chen, J.-M., Cooper, D. N., Chuzhanova, N., Férec, C., and Patrinos, G. P. 2007. Gene conversion: mechanisms,  
450 evolution and human disease. *Nature Reviews Genetics*, 8(10): 762–775.
- 451 Comeron, J. M., Ratnappan, R., and Bailin, S. 2012. The Many Landscapes of Recombination in *Drosophila*  
452 *melanogaster*. *PLOS Genetics*, 8(10): e1002905. Publisher: Public Library of Science.
- 453 Corcoran, P., Gossmann, T. I., Barton, H. J., Great Tit HapMap Consortium, Slate, J., and Zeng, K. 2017. Determi-  
454 nants of the Efficacy of Natural Selection on Coding and Noncoding Variability in Two Passerine Species. *Genome*  
455 *Biol Evol*, 9(11): 2987–3007.
- 456 de Procé, S. M., Zeng, K., Betancourt, A. J., and Charlesworth, B. 2012. Selection on codon usage and base  
457 composition in *Drosophila americana*. *Biology Letters*, 8(1): 82–85. Publisher: Royal Society.
- 458 Duret, L. and Arndt, P. F. 2008. The Impact of Recombination on Nucleotide Substitutions in the Human Genome.  
459 *PLOS Genetics*, 4(5): e1000071.
- 460 Duret, L. and Galtier, N. 2009. Biased Gene Conversion and the Evolution of Mammalian Genomic Landscapes.  
461 *Annual Review of Genomics and Human Genetics*, 10(1): 285–311.
- 462 Ellegren, H., Smeds, L., Burri, R., Olason, P. I., Backström, N., Kawakami, T., Künstner, A., Mäkinen, H.,  
463 Nadachowska-Brzyska, K., Qvarnström, A., Uebbing, S., and Wolf, J. B. W. 2012. The genomic landscape of  
464 species divergence in *Ficedula flycatchers*. *Nature*, 491(7426): 756–760.
- 465 Eyre-Walker, A. and Hurst, L. D. 2001. The evolution of isochores. *Nature Reviews Genetics*, 2(7): 549.
- 466 Eyre-Walker, A., Woolfit, M., and Phelps, T. 2006. The distribution of fitness effects of new deleterious amino acid  
467 mutations in humans. *Genetics*, 173(2): 891–900.
- 468 Galtier, N. and Duret, L. 2007. Adaptation or biased gene conversion? Extending the null hypothesis of molecular  
469 evolution. *Trends in Genetics*, 23(6): 273–277.
- 470 Galtier, N., Roux, C., Rousselle, M., Romiguier, J., Figuet, E., Glémin, S., Bierne, N., and Duret, L. 2018. Codon  
471 Usage Bias in Animals: Disentangling the Effects of Natural Selection, Effective Population Size, and GC-Biased  
472 Gene Conversion. *Molecular Biology and Evolution*, 35(5): 1092–1103.
- 473 Garrison, E. and Marth, G. 2012. Haplotype-based variant detection from short-read sequencing. *arXiv:1207.3907*  
474 *[q-bio]*.
- 475 Glémin, S., Arndt, P. F., Messer, P. W., Petrov, D., Galtier, N., and Duret, L. 2015. Quantification of GC-biased  
476 gene conversion in the human genome. *Genome Research*, 25(8): 1215–1228.
- 477 Gossmann, T. I., Bockwoldt, M., Diringer, L., Schwarz, F., and Schumann, V.-F. 2018. Evidence for Strong Fixation  
478 Bias at 4-fold Degenerate Sites Across Genes in the Great Tit Genome. *Frontiers in Ecology and Evolution*, 6.

- 479 Groemping, U. 2006. Relative Importance for Linear Regression in R: The Package relaimpo. *Journal of Statistical*  
480 *Software*, 17(1): 1–27. Number: 1.
- 481 Gutz, H. and Leslie, J. F. 1976. Gene Conversion: A Hitherto Overlooked Parameter in Population Genetics. *Genetics*,  
482 83(4): 861–866.
- 483 Haddrill, P. R. and Charlesworth, B. 2008. Non-neutral processes drive the nucleotide composition of non-coding  
484 sequences in *Drosophila*. *Biology Letters*, 4(4): 438–441.
- 485 Haddrill, P. R., Bachtrog, D., and Andolfatto, P. 2008. Positive and Negative Selection on Noncoding DNA in  
486 *Drosophila simulans*. *Molecular Biology and Evolution*, 25(9): 1825–1834.
- 487 Halligan, D. L. and Keightley, P. D. 2006. Ubiquitous selective constraints in the *Drosophila* genome revealed by a  
488 genome-wide interspecies comparison. *Genome Research*, 16(7): 875–884.
- 489 Hansson, B., Ljungqvist, M., Dawson, D. A., Mueller, J. C., Olano-Marin, J., Ellegren, H., and Nilsson, J.-A. 2010.  
490 Avian genome evolution: insights from a linkage map of the blue tit (*Cyanistes caeruleus*). *Heredity*, 104(1): 67–78.
- 491 Harris, R. S. 2007. Improved pairwise alignment of genomic DNA. *Ph.D. Thesis, The Pennsylvania State University*.
- 492 Hayes, K., Barton, H. J., and Zeng, K. 2020. A study of faster-Z evolution in the great tit (*Parus major*). *Genome*  
493 *Biology and Evolution*.
- 494 Hillier, L. W., Miller, W., Birney, E., Warren, W., Hardison, R. C., Ponting, C. P., Bork, P., Burt, D. W., Groenen,  
495 M. A., Delany, M. E., and others 2004. Sequence and comparative analysis of the chicken genome provide unique  
496 perspectives on vertebrate evolution. *Nature*, 432(7018): 695–716.
- 497 Hodgkinson, A. and Eyre-Walker, A. 2011. Variation in the mutation rate across mammalian genomes. *Nat. Rev.*  
498 *Genet.*, 12(11): 756–766.
- 499 Hwang, D. G. and Green, P. 2004. Bayesian Markov chain Monte Carlo sequence analysis reveals varying neutral  
500 substitution patterns in mammalian evolution. *Proceedings of the National Academy of Sciences of the United*  
501 *States of America*, 101(39): 13994–14001.
- 502 Jackson, B. C., Campos, J. L., Haddrill, P. R., Charlesworth, B., and Zeng, K. 2017. Variation in the Intensity of  
503 Selection on Codon Bias over Time Causes Contrasting Patterns of Base Composition Evolution in *Drosophila*.  
504 *Genome Biology and Evolution*, 9(1): 102–123.
- 505 Kent, W. J., Baertsch, R., Hinrichs, A., Miller, W., and Haussler, D. 2003. Evolution’s cauldron: Duplication,  
506 deletion, and rearrangement in the mouse and human genomes. *Proceedings of the National Academy of Sciences*,  
507 100(20): 11484–11489.

- 508 Kunstner, A., Nabholz, B., and Ellegren, H. 2011. Significant Selective Constraint at 4-Fold Degenerate Sites in  
509 the Avian Genome and Its Consequence for Detection of Positive Selection. *Genome Biology and Evolution*, 3(0):  
510 1381–1389.
- 511 Laine, V. N., Gossmann, T. I., Schachtschneider, K. M., Garroway, C. J., Madsen, O., Verhoeven, K. J. F., de Jager,  
512 V., Megens, H.-J., Warren, W. C., Minx, P., Crooijmans, R. P. M. A., Corcoran, P., Great Tit HapMap Consortium,  
513 Sheldon, B. C., Slate, J., Zeng, K., van Oers, K., Visser, M. E., and Groenen, M. A. M. 2016. Evolutionary signals  
514 of selection on cognition from the great tit genome and methylome. *Nat Commun*, 7: 10474.
- 515 Liu, H., Huang, J., Sun, X., Li, J., Hu, Y., Yu, L., Liti, G., Tian, D., Hurst, L. D., and Yang, S. 2018. Tetrad analysis  
516 in plants and fungi finds large differences in gene conversion rates but no GC bias. *Nature Ecology & Evolution*,  
517 2(1): 164–173.
- 518 Long, H., Sung, W., Kucukyildirim, S., Williams, E., Miller, S. F., Guo, W., Patterson, C., Gregory, C., Strauss, C.,  
519 Stone, C., Berne, C., Kysela, D., Shoemaker, W. R., Muscarella, M. E., Luo, H., Lennon, J. T., Brun, Y. V., and  
520 Lynch, M. 2018. Evolutionary determinants of genome-wide nucleotide composition. *Nature Ecology & Evolution*,  
521 2(2): 237–240.
- 522 Matsumoto, T., Akashi, H., and Yang, Z. 2015. Evaluation of Ancestral Sequence Reconstruction Methods to Infer  
523 Nonstationary Patterns of Nucleotide Substitution. *Genetics*, 200(3): 873–890.
- 524 McVean, G. a. T. and Charlesworth, B. 1999. A population genetic model for the evolution of synonymous codon  
525 usage: patterns and predictions. *Genetics Research*, 74(2): 145–158.
- 526 Mugal, C. F., Arndt, P. F., and Ellegren, H. 2013. Twisted signatures of GC-biased gene conversion embedded in an  
527 evolutionary stable karyotype. *Molecular Biology and Evolution*, 30(7): 1700–1712.
- 528 Muyle, A., Serres-Giardi, L., Ressayre, A., Escobar, J., and Glémin, S. 2011. GC-Biased Gene Conversion and  
529 Selection Affect GC Content in the *Oryza* Genus (rice). *Molecular Biology and Evolution*, 28(9): 2695–2706.
- 530 Nagylaki, T. 1983. Evolution of a finite population under gene conversion. *Proceedings of the National Academy of*  
531 *Sciences*, 80(20): 6278–6281.
- 532 Parvanov, E. D., Petkov, P. M., and Paigen, K. 2010. Prdm9 Controls Activation of Mammalian Recombination  
533 Hotspots. *Science*, 327(5967): 835–835.
- 534 R Core Team 2015. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Com-  
535 puting, Vienna, Austria.
- 536 Rajic, Z. A., Jankovic, G. M., Vidovic, A., Milic, N. M., Skoric, D., Pavlovic, M., and Lazarevic, V. 2005. Size of the  
537 protein-coding genome and rate of molecular evolution. *Journal of Human Genetics*, 50(5): 217–229.

- 538 Ratnakumar, A., Mousset, S., Glémin, S., Berglund, J., Galtier, N., Duret, L., and Webster, M. T. 2010. Detecting  
539 positive selection within genomes: the problem of biased gene conversion. *Philosophical Transactions of the Royal  
540 Society B: Biological Sciences*, 365(1552): 2571–2580.
- 541 Rousselle, M., Laverre, A., Figuet, E., Nabholz, B., and Galtier, N. 2019. Influence of Recombination and GC-biased  
542 Gene Conversion on the Adaptive and Nonadaptive Substitution Rate in Mammals versus Birds. *Molecular Biology  
543 and Evolution*, 36(3): 458–471.
- 544 Ségurel, L., Wyman, M. J., and Przeworski, M. 2014. Determinants of Mutation Rate Variation in the Human  
545 Germline. *Annual Review of Genomics and Human Genetics*, 15(1): 47–70.
- 546 Singhal, S., Leffler, E. M., Sannareddy, K., Turner, I., Venn, O., Hooper, D. M., Strand, A. I., Li, Q., Raney, B.,  
547 Balakrishnan, C. N., Griffith, S. C., McVean, G., and Przeworski, M. 2015. Stable recombination hotspots in birds.  
548 *Science*, 350(6263): 928–32.
- 549 Smeds, L., Mugal, C. F., Qvarnström, A., and Ellegren, H. 2016. High-Resolution Mapping of Crossover and Non-  
550 crossover Recombination Events by Whole-Genome Re-sequencing of an Avian Pedigree. *PLOS Genetics*, 12(5):  
551 e1006044.
- 552 Stapley, J., Birkhead, T. R., Burke, T., and Slate, J. 2008. A Linkage Map of the Zebra Finch *Taeniopygia guttata*  
553 Provides New Insights Into Avian Genome Evolution. *Genetics*, 179(1): 651–667.
- 554 Stapley, J., Feulner, P. G. D., Johnston, S. E., Santure, A. W., and Smadja, C. M. 2017. Variation in recombination  
555 frequency and distribution across eukaryotes: patterns and processes. *Phil. Trans. R. Soc. B*, 372(1736): 20160455.
- 556 Van der Auwera, G. A., Carneiro, M. O., Hartl, C., Poplin, R., Del Angel, G., Levy-Moonshine, A., Jordan, T., Shakir,  
557 K., Roazen, D., Thibault, J., Banks, E., Garimella, K. V., Altshuler, D., Gabriel, S., and DePristo, M. A. 2013.  
558 From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Current  
559 Protocols in Bioinformatics / Editorial Board, Andreas D. Baxeavanis ... [et Al.]*, 43: 11.10.1–33.
- 560 van Oers, K., Santure, A. W., De Cauwer, I., van Bers, N. E., Crooijmans, R. P., Sheldon, B. C., Visser, M. E., Slate,  
561 J., and Groenen, M. A. 2014. Replicated high-density genetic maps of two great tit populations reveal fine-scale  
562 genomic departures from sex-equal recombination rates. *Heredity*, 112(3): 307–316.
- 563 Wallberg, A., Glémin, S., and Webster, M. T. 2015. Extreme Recombination Frequencies Shape Genome Variation  
564 and Evolution in the Honeybee, *Apis mellifera*. *PLOS Genetics*, 11(4): e1005189.
- 565 Warren, W. C., Clayton, D. F., Ellegren, H., Arnold, A. P., Hillier, L. W., Künstner, A., Searle, S., White, S., Vilella,  
566 A. J., Fairley, S., Heger, A., Kong, L., Ponting, C. P., Jarvis, E. D., Mello, C. V., Minx, P., Lovell, P., Velho, T.  
567 A. F., Ferris, M., Balakrishnan, C. N., Sinha, S., Blatti, C., London, S. E., Li, Y., Lin, Y.-C., George, J., Sweedler,  
568 J., Southey, B., Gunaratne, P., Watson, M., Nam, K., Backström, N., Smeds, L., Nabholz, B., Itoh, Y., Whitney,  
569 O., Pfenning, A. R., Howard, J., Völker, M., Skinner, B. M., Griffin, D. K., Ye, L., McLaren, W. M., Flicek,

- 570 P., Quesada, V., Velasco, G., Lopez-Otin, C., Puente, X. S., Olender, T., Lancet, D., Smit, A. F. A., Hubley, R.,  
571 Konkel, M. K., Walker, J. A., Batzer, M. A., Gu, W., Pollock, D. D., Chen, L., Cheng, Z., Eichler, E. E., Stapley, J.,  
572 Slate, J., Ekblom, R., Birkhead, T., Burke, T., Burt, D., Scharff, C., Adam, I., Richard, H., Sultan, M., Soldatov,  
573 A., Lehrach, H., Edwards, S. V., Yang, S.-P., Li, X., Graves, T., Fulton, L., Nelson, J., Chinwalla, A., Hou, S.,  
574 Mardis, E. R., and Wilson, R. K. 2010. The genome of a songbird. *Nature*, 464(7289): 757–762.
- 575 Weber, C. C., Boussau, B., Romiguier, J., Jarvis, E. D., and Ellegren, H. 2014. Evidence for GC-biased gene conversion  
576 as a driver of between-lineage differences in avian base composition. *Genome Biology*, 15(12): 549.
- 577 Williams, A. L., Genovese, G., Dyer, T., Altemose, N., Truax, K., Jun, G., Patterson, N., Myers, S. R., Curran, J. E.,  
578 Duggirala, R., Blangero, J., Reich, D., and Przeworski, M. 2015. Non-crossover gene conversions show strong GC  
579 bias and unexpected clustering in humans. *eLife*, 4: e04637.
- 580 Yang, Z. 2007. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology and Evolution*, 24(8):  
581 1586–1591.
- 582 Zeng, K., Jackson, B. C., and Barton, H. J. 2019. Methods for Estimating Demography and Detecting Between-Locus  
583 Differences in the Effective Population Size and Mutation Rate. *Molecular Biology and Evolution*, 36(2): 423–433.
- 584 Zhang, G., Li, C., Li, Q., Li, B., Larkin, D. M., Lee, C., Storz, J. F., Antunes, A., Greenwold, M. J., Meredith,  
585 R. W., Ödeen, A., Cui, J., Zhou, Q., Xu, L., Pan, H., Wang, Z., Jin, L., Zhang, P., Hu, H., Yang, W., Hu, J.,  
586 Xiao, J., Yang, Z., Liu, Y., Xie, Q., Yu, H., Lian, J., Wen, P., Zhang, F., Li, H., Zeng, Y., Xiong, Z., Liu, S.,  
587 Zhou, L., Huang, Z., An, N., Wang, J., Zheng, Q., Xiong, Y., Wang, G., Wang, B., Wang, J., Fan, Y., da Fonseca,  
588 R. R., Alfaro-Núñez, A., Schubert, M., Orlando, L., Mourier, T., Howard, J. T., Ganapathy, G., Pfenning, A.,  
589 Whitney, O., Rivas, M. V., Hara, E., Smith, J., Farré, M., Narayan, J., Slavov, G., Romanov, M. N., Borges, R.,  
590 Machado, J. P., Khan, I., Springer, M. S., Gatesy, J., Hoffmann, F. G., Opazo, J. C., Håstad, O., Sawyer, R. H.,  
591 Kim, H., Kim, K.-W., Kim, H. J., Cho, S., Li, N., Huang, Y., Bruford, M. W., Zhan, X., Dixon, A., Bertelsen,  
592 M. F., Derryberry, E., Warren, W., Wilson, R. K., Li, S., Ray, D. A., Green, R. E., O'Brien, S. J., Griffin, D.,  
593 Johnson, W. E., Haussler, D., Ryder, O. A., Willerslev, E., Graves, G. R., Alström, P., Fjeldså, J., Mindell, D. P.,  
594 Edwards, S. V., Braun, E. L., Rahbek, C., Burt, D. W., Houde, P., Zhang, Y., Yang, H., Wang, J., Avian Genome  
595 Consortium, Jarvis, E. D., Gilbert, M. T. P., and Wang, J. 2014. Comparative genomics reveals insights into avian  
596 genome evolution and adaptation. *Science (New York, N.Y.)*, 346(6215): 1311–1320.