1 # Selective sweeps influence diversity over

2 # large regions of the mouse genome

3

4 Tom R. Booker[1,2], Benjamin C. Jackson[3], Rory J. Craig[3], Brian Charlesworth[3], Peter D. Keightley[3]

5

6 1.     Department of Zoology, University of British Columbia, Vancouver, British Columbia,

7 Canada.

8 2.     Department of Biological Sciences, University of Calgary, Calgary, Alberta, Canada.

9 3.     Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK. EH9 3FL

10

11 Correspondence: booker@zoology.ubc.ca

12

13

# Abstract

To what extent do substitutions in protein-coding versus gene-regulatory regions contribute to fitness change over time? Answering this question requires estimates of the extent of selection acting on beneficial mutations in the two classes of sites. New mutations that have advantageous or deleterious fitness effects can induce selective sweeps and background selection, respectively, causing variation in the level of neutral genetic diversity along the genome. In this study, we analyse the profiles of genetic variability around protein-coding and regulatory elements in the genomes of wild mice to estimate the parameters of positive selection. We find patterns of diversity consistent with the effects of selection at linked sites, which are similar across mouse taxa, despite differences in effective population size and demographic history. By fitting a model that combines the effects of selective sweeps and background selection, we estimate the strength of positive selection and the frequency of advantageous mutations. We find that strong positive selection is required to explain variation in genetic diversity across the murid genome. In particular, we estimate that beneficial mutations in protein-coding regions have stronger effects on fitness than do mutations in gene-regulatory regions, but that mutations in gene-regulatory regions are more common. Overall though, our parameter estimates suggest that the cumulative fitness changes brought about by beneficial mutations in protein-coding may be greater than those in gene-regulatory elements.

# Introductions

Understanding the relative contributions of protein-coding and gene regulatory variation to adaptation is a long-standing goal of evolutionary biology. Molecular changes in protein-coding and gene regulatory regions contribute to evolution, but in classic essays King and Wilson (1975) and Carroll (2005) argued that changes in gene expression may dominate adaptive evolution. King and Wilson (1975) reasoned that since nucleotide identity between human and chimpanzee proteins is around 99%, there are too few protein sequence difference between the species, implying that changes in gene regulation are probably required to explain the many phenotypic differences between the species. Carroll's (2005) argument highlighted the idea that molecular changes in the gene regulatory apparatus may have smaller pleiotropic effects than those in protein-coding regions, so that changes in gene expression may dominate adaptive evolution. However, Hoekstra and Coyne (2007) attempted to refute these arguments, maintaining that there is insufficient evidence to decide whether adaptation is primarily driven by changes in protein sequences or gene regulatory elements. For example, a 1% difference in protein sequence between humans and chimpanzees could still result in a very large number of phenotypic differences. However, the contribution of individual variants to additive genetic variance for a trait is expected to be proportional to the square of their phenotypic effect sizes, assuming semi-dominance (Fisher 1918; Falconer and Mackay 1996). Without an understanding of the frequencies of new mutations, their effect sizes, and the strength of selection acting on them, the question of the contribution of molecular evolution in different genomic elements to adaptation will remain intractable.

Information on the strength of selection acting on beneficial mutations and the rates at which they occur can be obtained by analysing patterns of neutral genetic diversity. Because genetically linked sites do not evolve independently, natural selection acting at a given site may leave signatures at linked sites that are informative about the strength and mode of selection. The effects of selection at linked sites on neutral genetic diversity depend on the frequency and strength of selected mutations and the rate of recombination (Charlesworth 2012; Hermisson and Pennings 2017; Stephan 2019). Several modes of selection at linked sites have been identified. Of specific relevance to this study are background selection (BGS), caused by the removal of deleterious mutations from a population, and selective sweeps, caused by the spread of advantageous variants. The classic footprint of a selective sweep is a trough in nucleotide diversity

3

67  at neutral sites surrounding an adaptive substitution. The reduction in nucleotide diversity caused

68  by a sweep is proportional to the ratio of the strength of selection acting on the causal mutation

69  to the local recombination rate (Barton 2000). Using such information, Wiehe and Stephan (1993)

70  developed a model of recurrent selective sweeps and used it to estimate the frequency and

71  strength of advantageous mutations in *Drosophila melanogaster*. They fitted their model of

72  sweeps to the relationship between recombination rate and nucleotide diversity for a number of

73  loci sampled across the *D. melanogaster* genome. At the time of their analysis, the theory of BGS

74  was in its infancy, and models combining the effects of BGS and sweeps had not been developed.

75  However, the effects of BGS are expected to be ubiquitous across the genome (McVicker et al.

76  2009; Comeron, 2014; Elyashiv et al. 2016; Pouyet et al. 2018), and conceptually similar studies to

77  Wiehe and Stephan (1993) have shown that controlling for BGS is important when parametrizing

78  sweep models (Kim and Stephan 2000; Comeron 2014; Elyashiv et al. 2016; Campos et al. 2017).

79

80  In *Drosophila*, there are reductions in average diversity around recent nonsynonymous

81  substitutions, which are greater than those observed around synonymous substitutions (Sattath et

82  al. 2011; Elyashiv et al. 2016). To investigate the causes of this difference, Elyashiv et al. (2016)

83  fitted a model of sweeps and BGS to genome-wide variation in genetic diversity in *D.*

84  *melanogaster* and found that a combination of BGS and selective sweeps provided a close fit to

85  the observed data. From the fit of their model to empirical data, Elyashiv et al. (2016) inferred a

86  distribution of fitness effects for advantageous mutations that included a class of very strongly

87  selected mutations and a more mildly beneficial class. In both mice and humans, however, there is

88  very little difference between the profiles of diversity around recent nonsynonymous and

89  synonymous substitutions (Hernandez et al. 2011; Halligan et al. 2013). In these species, dips in

90  average nucleotide diversity have been observed in genomic regions surrounding whole functional

91  elements, such as protein-coding exons or conserved non-coding elements, which may reflect the

92  cumulative effects of recurrent selective sweeps and BGS (Hernandez et al. 2011; Halligan et al.

93  2013; Booker and Keightley 2018)

94

95  Natural populations of mice in the genus *Mus* are excellent material for the study of adaptive

96  evolution in different regions of the mammalian genome. Their populations are very large

97  compared to other mammals (Leffler et al. 2012), so there is likely to be more power for

98  population genetic analyses to differentiate between the evolutionary processes that affect

99  genetic variability. Previous studies in mice have established that both protein-coding genes and

100 regions putatively involved in gene regulation have an excess of sequence differences from sister

101 taxa compared to that expected under a model of purifying selection, suggesting widespread

102 adaptive molecular evolution (Halligan et al. 2010, 2013). Halligan et al. (2013) analysed a sample

103 of *Mus musculus castaneus* individuals and estimated that there have been around 1.3 million and

104 0.38 million positively selected regulatory and nonsynonymous changes, respectively, over the

105 period since this subspecies began to diverge from rats. At face value, this finding suggests that

106 changes in gene regulation may dominate adaptive evolution in mice. However, Halligan et al.

107 (2013) also showed that there are much larger reductions in neutral diversity surrounding protein-

108 coding exons than around gene regulatory elements, and that BGS could not fully explain these

109 observations (Halligan et al. 2013; Booker and Keightley 2018). Halligan et al. (2013) concluded

110 that this difference in neutral diversity may reflect differences in the strength of positive selection

111 acting on the different classes of sites.

112

113 Building on Halligan et al. (2013), we have sought to tease apart the contributions of BGS and

114 sweeps to the patterns of nucleotide diversity observed in the Eastern house mouse *M. m.*

115 *castaneus* (Booker and Keightley 2018). We inferred the distribution of fitness effects (DFE) for

116 deleterious and advantageous mutations occurring in protein-coding genes and gene regulatory

117 elements, by analysing the frequency distribution of derived allele frequencies (the unfolded site

118 frequency spectrum, uSFS). Based on analysis of the uSFS, we found that a model of positive

119 selection was insufficient to explain the troughs in nucleotide diversity around protein-coding

120 exons or conserved non-coding elements (CNEs). However, we found that infrequent, strongly

121 beneficial mutations, which have negligible effect on the uSFS, potentially could do so (Booker and

122 Keightley 2018). This is because infrequent, strongly advantageous mutations may substantially

123 influence diversity at linked sites, while making very little contribution to the uSFS. We concluded

124 that the parameters of positive selection are very difficult to accurately estimate from the uSFS

125 alone (Booker 2020). To understand the relative strengths of selection acting on protein-coding

126 versus gene regulatory regions, the analysis of a model of selective sweeps fitted to patterns of

127 neutral genetic variability may be more powerful.

128

129 In this study, we examine the reductions in nucleotide diversity surrounding protein-coding exons

130 and conserved non-coding elements in wild mice, and attempt to tease apart the modes of

131 selection operating on the two different elements. We fitted a model of selective sweeps to the

132 patterns observed in *M. m. castaneus,* while correcting for the confounding effects of BGS. Our

5

133     analysis provides evidence that the strength of selection acting on beneficial mutations in protein-

134     coding exons is far greater than that acting on conserved non-coding elements. Using a simple

135     model of the fitness change brought about by positive selection, we find that selection on protein-

136     coding regions may contribute more to fitness change, despite positive selection occurring more

137     frequently in regulatory regions of the genome. We then compared patterns of putatively neutral

138     diversity among the principal subspecies of *Mus musculus* and their sister species *Mus spretus*. We

139     find that the profiles of nucleotide diversity and the inferred distributions of fitness effects among

140     each group are similar, suggesting that the contributions of positive selection to protein-coding

141     and regulatory change are similar in the different mouse taxa. Note that our goal in this study is

142     not to identify the individual loci that selection has recently acted on; for a recent study

143     identifying the targets of recent selection in wild mice see (Lawal et al. 2021).

# Results and Discussion

*Profiles of genetic diversity around protein-coding exons and conserved non-coding elements in multiple mouse lineages*

If different mouse lineages are subject to similar selection pressures, we might expect that they exhibit similar patterns of diversity across their genomes. We thus compared patterns of genetic diversity in populations of the house mouse *Mus musculus* and the sister species *Mus spretus.* We analysed data previously reported by Halligan et al. (2013) and Harr et al. (2016) for the two mouse species, *M. musculus* and *M. spretus.* For *M. musculus*, we analysed samples from the three sub-species, *M. m. castaneus*, *M. m. domesticus* and *M. m. musculus*. The *M. m. castaneus* individuals (*n* = 10) were from Himachal Pradesh, India. For *M. m. domesticus,* populations were sampled in France (*n* = 8), Germany (*n* = 8) and Iran (*n* = 8). In the case of *M. m. musculus,* populations were sampled in Afghanistan (*n* = 6), the Czech Republic (*n* = 8) and Kazakhstan (*n* = 8). The *M. spretus* individuals were sampled in Spain (*n* = 8). We refer to the different sub-populations of *M. m. domesticus* and *M. m. musculus* by the countries where the individuals were sampled.

We identified conserved non-coding elements (CNEs) in murid rodents using a 40-way alignment of placental mammals by means of the *phastCons* approach (Siepel et al. 2005). Following Williamson et al. (2014), the genomes of *M. musculus* and other rodents were masked in the alignment to limit ascertainment bias affecting elements that have recently diverged in the rodent lineage. CNEs identified using *phastCons* overlap with features such as promoters and enhancers (Lindblad-Toh et al. 2011), and thus are likely to have roles in the regulation of gene expression.

For each of the mouse taxa, we examined putatively neutral nucleotide diversity ($\pi$) in genomic regions surrounding protein-coding exons and CNEs using the methods described by Halligan et al. (2013). Briefly, polymorphism data were extracted in genomic windows surrounding protein-coding exons and CNEs. We masked any putatively functional sites from analysis windows; these included the exons (including UTRs) of genes annotated in the *M. musculus* genome by ENSEMBL in release 93 (Howe et al. 2021) and CNEs. For each analysis window, we calculated the genetic map distance between the centre of the window and the focal functional element, assuming either the pedigree-based recombination map for *M. musculus* constructed by Cox et al. (2009) or a recombination map estimated using linkage disequilibrium (LD) in the *M. m. castaneus* genome

7

176   (Appendix). We excluded analysis windows that had a scaled genetic distance of $4N_er$ < 1, because

177   downstream analyses assume that sites are not tightly linked. All remaining analysis windows

178   were collated into genetic distance bins. The average number of pairwise differences (i.e.

179   nucleotide diversity) and nucleotide divergence from *Rattus rattus* were calculated for each

180   genetic distance bin. For all analyses, we only examined the autosomes. Downstream analyses

181   were sensitive to the assumption of a single mutation rate, so we excluded hypermutable CpG-

182   prone sites from our analyses, identified as sites that were preceded by a C or succeeded by a G in

183   the 3' to 5' direction in the reference genome (see Materials and Methods).

184

185   The choice of recombination map had a substantial effect on the profiles of average nucleotide

186   diversity observed around protein-coding exons and CNEs (Figures S1, S2). When assuming the

187   pedigree-based Cox map, we found that nucleotide diversity was slightly higher in the immediate

188   flanks of both exons and CNEs (with distances calculated using the estimated recombination

189   frequency), and lower in regions far from functional elements, compared to results with the LD-

190   based map (Figure S1, S2). These differences are consistent with the possibility that the Cox map,

191   which was constructed with a far smaller number of markers than the LD-based map, does not

192   fully capture genomic regions that have either unusually low or high recombination rates. A

193   possible consequence of this would be that analysis windows at intermediate distances from

194   functional elements may appear to be more tightly linked to those elements than they actually

195   are. This is supported by differences in numbers of sites falling into various genetic distance bins

196   between the pedigree-based and LD-based recombination maps (Figures S1 and S2). However,

197   both selective sweeps and BGS can induce LD, and may thus downwardly bias recombination rate

198   estimates obtained using LD-based approaches (Clark et al., 2010). For this reason, we focus on

199   results obtained assuming the pedigree-based Cox map for the remainder of the paper. We

200   present parallel analyses, which assume the LD-based map, in the supplement and describe

201   differences between the respective conclusions in the Discussion.
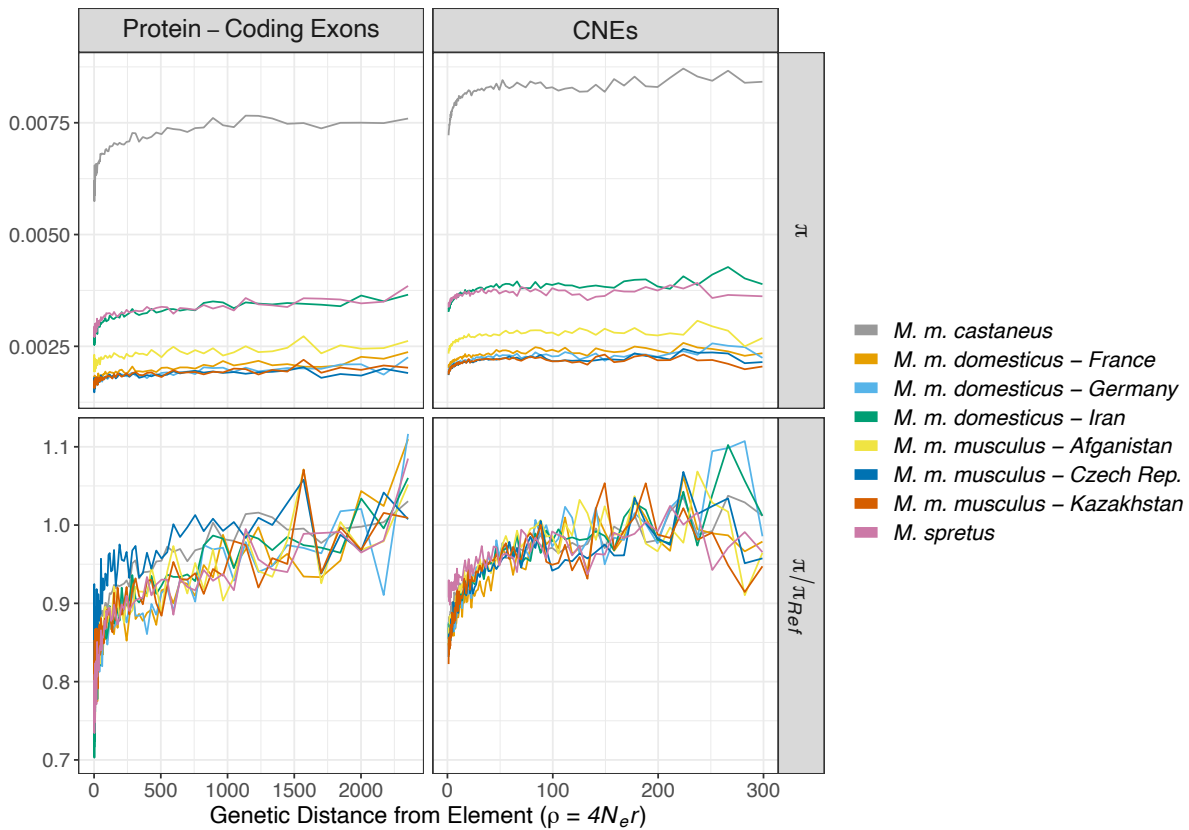
202

203   All mouse taxa exhibited dips in nucleotide diversity around protein-coding exons and CNEs

204   (Figure 1). To quantify the relative reduction in diversity for each taxa, we calculated $\pi/\pi_{Ref}$, the

205   ratio of $\pi$ to the average $\pi$ at distances greater than $4N_er$ = 1,500 and less than $4N_er$ = 2,500 for

206   exons, and distances greater than $4N_er$ = 150 and less than $4N_er$ = 250 for CNEs. The distances for

207   determining $\pi_{Ref}$ were chosen based on where $\pi$ began to flatten off with increasing distance from

208   functional elements. Despite the existence of large differences in genome-wide diversity between

209 the taxa, troughs in $\pi/\pi_{Ref}$ around exons and CNEs were very similar among mouse lineages

210 (Figures 1, S1, S2). Nucleotide diversity was reduced by 20-30% and 10-20% around protein-coding

211 exons and CNEs, respectively (Figure 1). The dips in diversity extended to genetic distances of up

212 to approximately $4N_e r = 1,000$ around exons, but only to $4N_e r = 100$ around CNEs (Figure 1).

213 Consistent with Halligan et al. (2013), we observed little reduction in between-species divergence

214 around the edges of protein-coding exons, suggesting that mutation rate variation is not a

215 substantial driver of the observed dips in diversity. However, in the immediate flanks of CNEs, we

216 observed a trough in divergence. This may be explained if the *phastCons* approach used to identify

217 CNEs did not readily identify weakly conserved sequences at the edge of more strongly conserved

218 blocks. This would imply that that some sites subject to purifying selection tightly linked to the

219 CNEs may have remained unannotated in our analysis. However, the troughs in nucleotide

220 divergence around CNEs were substantially narrower than the corresponding troughs in diversity.

221 This implies that reduced mutation rates or constrained sites may account for part of the diversity

222 drop around CNEs, but do not explain all of it (Figure S1, S2).

223

224 An important caveat concerning the above analysis is that the mouse taxa thatare subject of the

225 analysis are very closely related, i.e. it has been estimated that the *M. musculus* sub-species

226 complex began to diverge around 350,000 years ago (Geraldes et al, 2011). Furthermore, Geraldes

227 et al (2011) found extensive shared nucleotide variation among the sub-species, and that the

228 average $F_{ST}$ among the members of the sub-species complex ranged from 0.43 to 0.72. Thus,

229 patterns of polymorphism identified in the species are likely to be highly non-independent, and

230 differences in $\pi$ between the groups presumably reflect fluctuations in population sizes.

231

**Figure 1** Nucleotide diversity ($\pi$) in regions surrounding protein-coding exons and CNEs in wild mice. Population-scaled recombination rates ($4N_er$) were calculated assuming the recombination map for *M. musculus* constructed by Cox et al. (2009). $\pi_{Ref}$ is the mean diversity calculated for sites far from functional elements.

Nucleotide polymorphism and divergence in wild mouse genomes

A first step for determining whether there was a consistent signal of natural selection across the mouse genomes, was to identify three classes of functional sites and two classes of putatively neutral sites as follows. For protein-coding gene orthologues between mouse and rat, we identified 0-fold degenerate nonsynonymous sites and UTRs, and used 4-fold degenerate sites as a neutral comparator. Protein-coding sites within the binding motifs of exonic splice enhancers including synonymous sites appear to be subject to purifying selection, implying that 4-fold sites located within them cannot be considered as neutral (Savisaar and Hurst 2018). We therefore excluded all synonymous and non-synonymous sites located within regions that matched such binding motifs. The total numbers of polymorphic and invariant sites that passed filtering for each

10

250 of the mouse taxa are detailed in Supplementary File 1. We identified sites in the upstream and

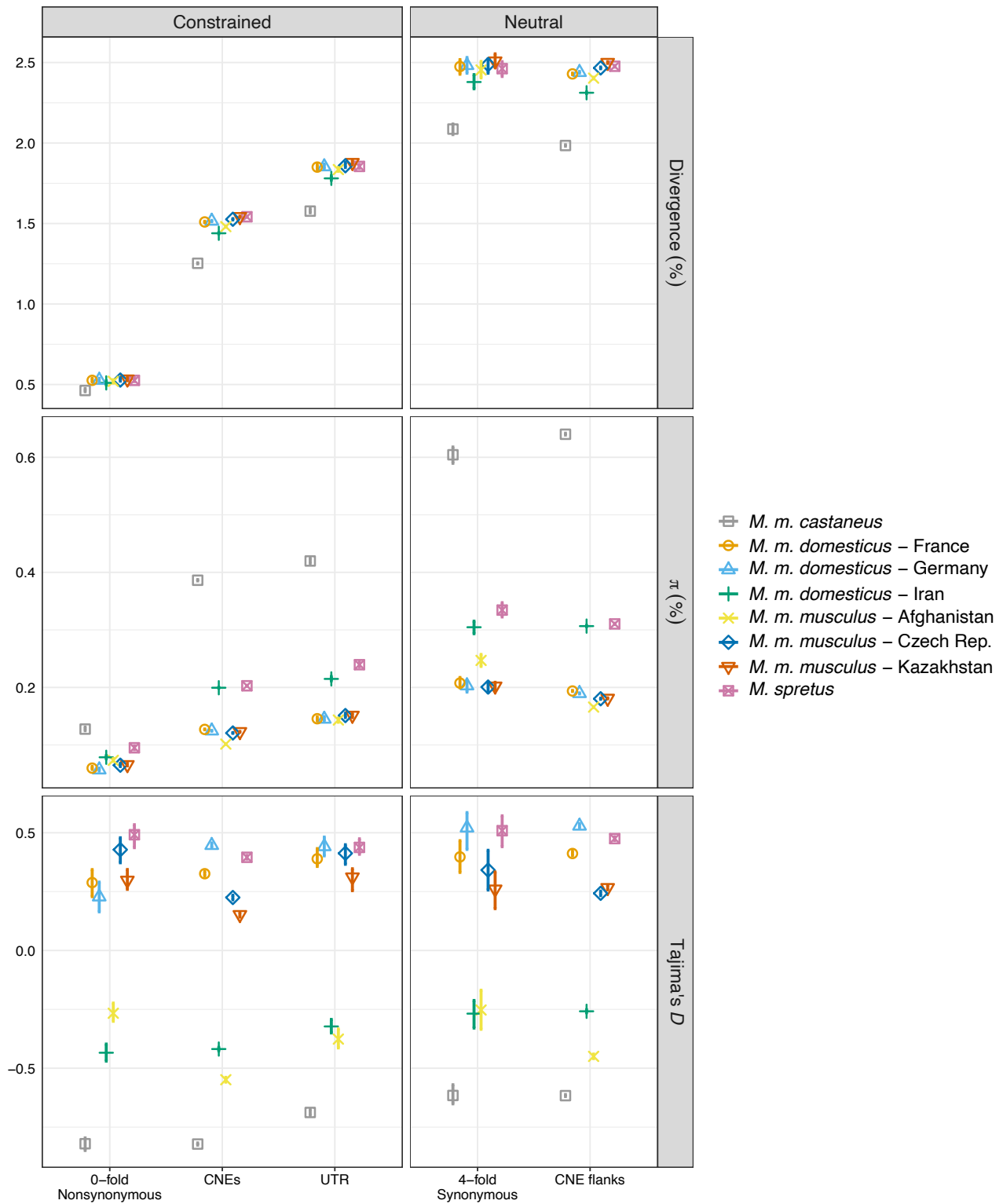251 downstream flanks of individual CNEs for use as neutral comparators (see Methods).

252

253 To determine whether there was a consistent signal of natural selection, we assessed nucleotide

254 diversity and lineage-specific divergence for the three classes of putatively functional sites (Figure

255 2).  In all cases, functional site diversity and divergence were lower than for their putatively

256 neutral counterparts, consistent with the action of purifying selection (Figure 2). Note that the *M.*

257 *m. castaneus* data have been analysed in this way before (Halligan et al. 2013; Booker and

258 Keightley 2018) and, as previously reported, *M. m. castaneus* had the highest nucleotide diversity

259 of all *Mus* taxa surveyed (Figure 2; Harr et al. 2016). However, nucleotide divergence reported is

260 the lineage specific divergence accumulated since the focal taxa began to diverge from *Mus*

261 *famulus*, so that divergence estimates for the various mouse taxa are highly non-independent

262 because of shared histories.

263

264 All populations had nonzero Tajima's *D* for putatively neutral sites, indicating the presence of

265 either non-equilibrium population dynamics or genome-wide effects of selection (Figure 2)*.* Mouse

266 populations from Western Europe and Kazakhstan exhibited positive Tajima's *D* for all classes of

267 sites (Figure 2), consistent with a recent history of admixture between different populations or

268 population bottlenecks (Charlesworth and Charlesworth 2010, pp.290-291). A population

269 structure analysis of the mice analysed in this study did not suggest strong admixture between the

270 sampled groups (Harr et al. 2016), but we cannot rule out the possibility of admixture with other

271 unsampled mouse populations. *M. m. castaneus* and populations sampled in Iran and Afghanistan

272 had strongly negative Tajima's *D* values, consistent with recent population expansion or a

273 genome-wide effect of recurrent selective sweeps (Charlesworth and Charlesworth 2010, pp.290,

274 414). Indeed, simulations modelling *D. melanogaster* populations have shown that recurrent,

275 strong selective sweeps can induce negative Tajima's *D* as large as -0.156 at synonymous sites

276 (Campos and Charlesworth 2019). It worth noting that Tajima's *D* is sensitive to the number of

277 individuals and nucleotides analysed (Simonsen et al. 1995), which vary among the mouse taxa, so

278 it is not straightforward to interpret differences in demographic history or strength of selection
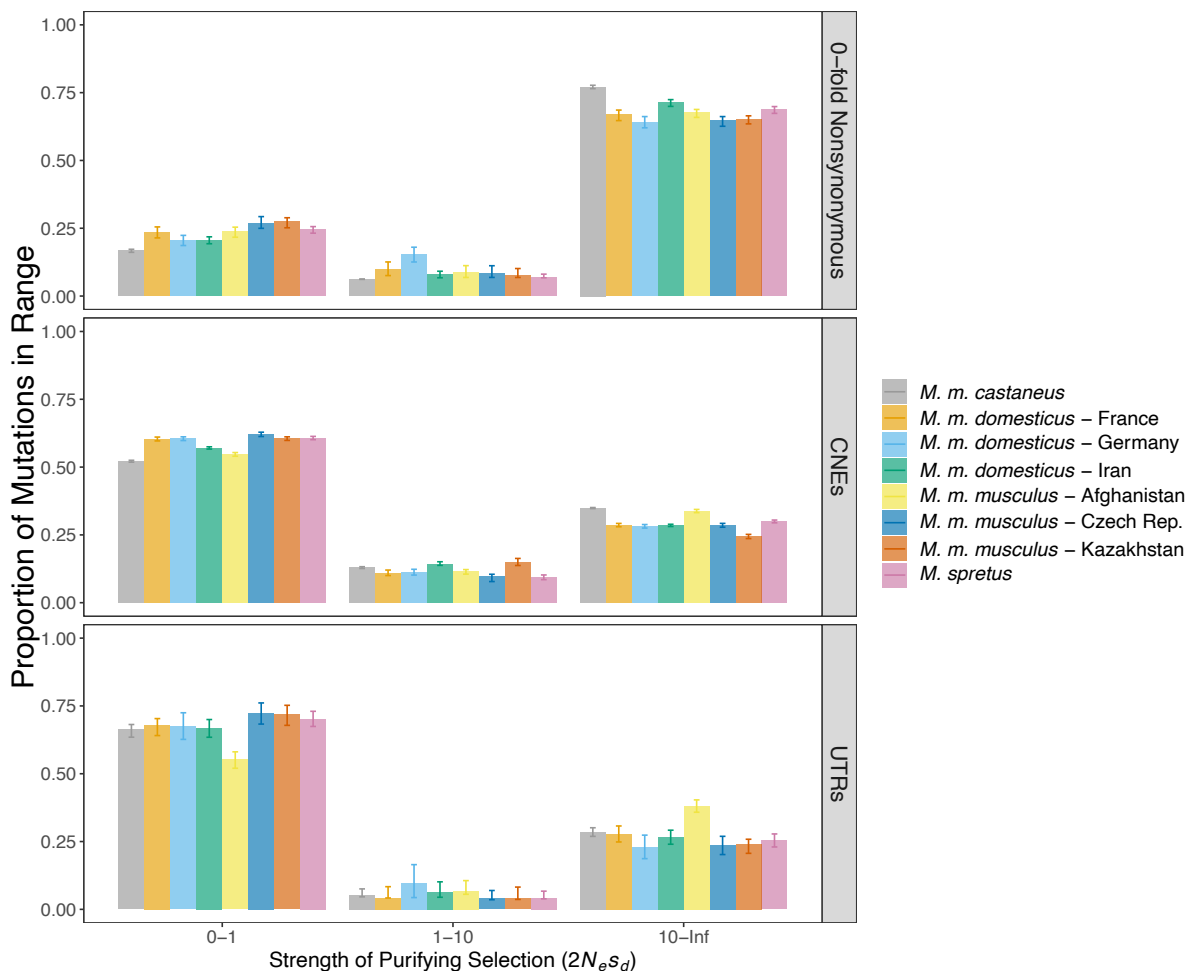
279 from these data.

280

**Figure 2** Population genetic summary statistics for three classes of putatively functional sites in the mouse genome and two putatively neutral comparators. Nucleotide diversity ($\pi$) and Tajima's *D* are also shown. Error bars represent 95% confidence intervals based on 100 bootstrap samples. Those not visible are shorter than the height of the points.

287     The distribution of fitness effects for deleterious mutations inferred from the uSFS

288     To parameterise a model of BGS, we estimated the distribution of fitness effects (DFE) for

289     deleterious mutations in each of the mouse taxa by fitting a model of mutation-selection-drift

290     balance to the unfolded site frequency spectrum (uSFS). The uSFS is a vector of 0, 1, 2, …, $k$ counts

291     of derived alleles, where $k$ is the number of haploid genomes sampled. Estimates of the DFE can

292     be obtained by contrasting the uSFS for a selected class of sites and a neutral comparator. Here,

293     we estimated the uSFS for the three classes of functional sites and their putatively neutral

294     comparator sequences. For each class of sites, we fitted a gamma distribution of deleterious

295     mutational effects using *polyDFE* (v2; Tataru and Bataillon 2019). Tataru et al. (2017) showed that

296     *polyDFE* provides robust estimates of the DFE for deleterious mutations based on the uSFS if a

297     discrete class of beneficial mutations is also inferred. While the inferred beneficial mutation

298     parameters are often spurious, including them seems to improve inference of the DFE for

299     deleterious mutations (Booker 2020). Finally, while the gamma distribution is an arbitrary choice

300     of model, and other probability distributions may give better fits to the data, it can capture the

301     important features of the DFE, even if the underlying distribution is multi-modal (Kousathanas and

302     Keightley 2013). Additionally, using the same probability distribution across taxa provides a

303     consistent framework for comparing molecular evolution in the different mouse groups.

304

305     The estimated DFEs were all highly leptokurtic and had similar estimated parameters across the

306     different taxa (Figure 3; Supplementary Table 2). Using *polyDFE,* the DFE is estimated in terms of

307     the scaled selection coefficient for deleterious mutations, $2N_es_d$, where $s_d$ is the reduction in

308     fitness experienced by an individual homozygous for the mutation (which is assumed to be semi-

309     dominant). Figure 3 shows the distribution of effects of deleterious mutational effects discretised

310     into three ranges; nearly neutral mutations with $2N_es_d < 1$, mildly deleterious mutations with $1 \leq$

311     $2N_es_d < 10$ and mutations with $2N_es_d \geq 10$. Consistent with previous studies, amino-acid changing

312     mutations were found to have the highest probability of having strongly deleterious effects

313     (Halligan et al. 2013) and non-coding elements (UTRs and CNEs) had higher fractions of nearly

314     neutral mutations (Figure 3). For 0-fold degenerate sites and CNEs, *M. m. castaneus* had the

315     smallest proportion of nearly neutral variants among the taxa. The DFE inferred for the *M. m.*

316     *musculus* sample from Afghanistan had the highest proportion of strongly deleterious mutations in

317     UTRs, but this may reflect sampling error, since there were only 6 individuals and the population

318     had among the lowest levels of nucleotide diversity (Figure 2).

319

13

**Figure 3** Graphical representation of the distribution of fitness effects of deleterious mutations for three classes of functional sites in wild mice. The figure shows the proportion of mutations falling into three ranges of effect size assuming a gamma DFE for each taxa and class of sites. Error bars indicate the 95% range based on 100 bootstrap replicates.

The contribution of background selection to patterns of diversity around functional elements

Using the inferred DFE parameters for deleterious mutations, we can estimate the contribution of BGS to reductions in nucleotide diversity across the mouse genome. Specifically, we used simulations modelling *M. m. castaneus* to estimate the contribution of BGS to troughs in diversity observed around functional elements (Figure S3). Our simulations incorporated recombination rate variation (assuming either the pedigree-based or LD-based recombination maps, see below), the distribution of exons, UTRs and CNEs in the mouse genome, and the distributions of fitness effects for deleterious mutations estimated for those elements. We estimated values of $\pi$ around both exons and CNEs from simulated data in the same manner as for the empirical data. Using the

336     simulation results, we estimated the reduction in diversity caused by background selection, *B,*

337     around functional elements for each genetic distance bin. We calculated $B = \pi/\pi_0$, where $\pi$ is the

338     nucleotide diversity observed in the simulation and $\pi_0$ is the neutral expectation. As we found

339     previously (Booker and Keightley 2018), BGS could not fully explain the reductions in diversity

340     observed around protein-coding exons or CNEs (Figure S3).

341

342     Inferences about the strength of BGS made under the assumption of constant population size be

343     misleading if there has been recent population size change. For example, a population bottleneck

344     may lead to the accumulation of weakly deleterious mutations if drift overwhelms selection. As

345     population size increases after a bottleneck, rapid purging of weakly deleterious mutations can

346     occur, leading to deviations from the expectations of standard models of BGS, which assume

347     constant population size (Torres et al. 2020; Johri et al. 2021). We have previously inferred a

348     model of demographic history for *M. m. castaneus*, which suggested that population size has

349     recently increased following a bottleneck (Booker and Keightley 2018). We performed an

350     additional set of simulations incorporating this demographic history, but found that the relative

351     reductions in diversity around both protein-coding exons and CNEs were very similar to those

352     observed under constant population size (Figure S4). Note that the trajectory of the demographic

353     history (bottleneck followed by recovery) we inferred may be an artefact of BGS (Ewing and

354     Jensen 2016; Johri et al. 2021). However, we proceeded with our analysis assuming estimates of *B*

355     for a constant population size, because the variations in *B* around exons and CNEs were very

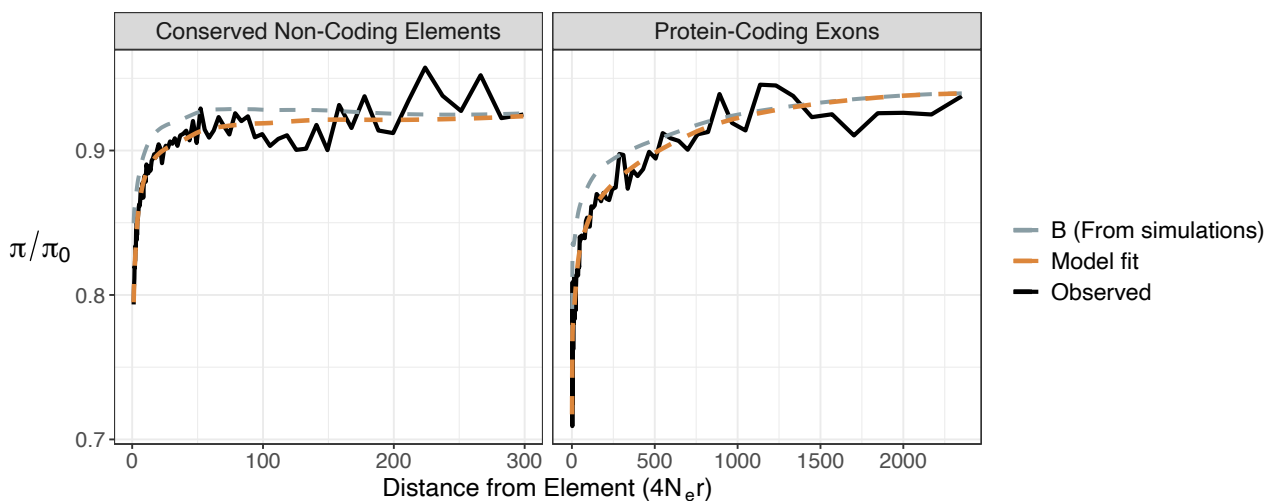356     similar with or without population size change.

357

358     *Parameters of beneficial mutations obtained from patterns of nucleotide diversity*

359     We estimated the parameters of beneficial mutations occurring in protein-coding and gene

360     regulatory regions by fitting a model that combines the effects of BGS and recurrent selective

361     sweeps to troughs in average nucleotide diversity around functional elements (see Materials and

362     Methods). The model quantifies the reduction in neutral diversity surrounding the average exon or

363     CNE, assuming that they are 150bp and 52bp long, respectively. A key parameter in the model is

364     $\pi_0 = 4N_e\mu$, the nucleotide diversity expected under neutrality in the absence of selection at linked

365     sites, where $N_e$ is the effective population size and $\mu$ is the mutation rate per basepair. Estimation

366     of $\pi_0$ is problematic, however, and $\pi_0$ may even be unobservable in empirical data, given the

367     ubiquity of selection at linked sites (Kern and Hahn 2018). In the empirical data, $\pi$ levelled off at

368    different values for protein-coding exons and CNEs (Figure 1, S1, S2). However, our simulations

369    predicted that $B$ should plateau at around 0.95 in genomic regions surrounding both protein-

370    coding exons and CNEs (Figure S3). $B$ was not predicted to plateau at 1.0 in our simulations,

371    because we modelled the distribution of all functional elements in the genome, so that a site may

372    be influenced by BGS generated by many surrounding elements. Our simulations did not model

373    sweeps, so simply dividing empirical $\pi$ by our estimated $B$ would give an underestimate of $\pi_0$,

374    because the reduction in diversity caused by positive selection was not included. When analysing

375    variation in $\pi$, we therefore assumed values of $\pi_0$ = 0.0081 and 0.0091 for protein-coding exons

376    and CNEs, respectively, to reflect the different levels at which diversity plateaued.

377

378    We proceeded to fit models combining BGS and selective sweeps to the troughs in diversity

379    around protein-coding exons and CNEs in *M. m. castaneus* assuming various models for the effects

380    of advantageous mutations. We estimated the strength of selection acting on new, semi-dominant

381    beneficial mutations as $\gamma_a = 4N_e s_a$, where $s_a$ is the increase in relative fitness experienced by

382    heterozygotes. We also estimated $p_a$, the proportion of new advantageous mutations in a

383    functional element. We found that a model with two classes of advantageous mutations gave a

384    better fit than a single class of mutations or an exponential distribution of effects (as judged by

385    AIC; Supplementary File 3). This result held regardless of the recombination map that was

386    assumed (Supplementary File 3).

387



388

389    **Figure 4** The reduction of scaled nucleotide diversity around protein-coding exons and CNEs in *M.*

390    *m. castaneus*, predicted by fitting a model combining the effects of background selection and

391    selective sweeps to the observed data. Genetic distances were calculated assuming the pedigree-

16

392    based recombination map constructed by Cox et al. (2009). The effect of background selection (*B*)

393    was estimated using simulations.

394

395    For both protein-coding exons and CNEs, we found that the best fitting model included a class of

396    strongly advantageous mutations and a class of more mildly beneficial mutations (Table 1). When

397    assuming the pedigree-based Cox map, we estimated the scaled fitness effects of the strongly

398    selected class ($\gamma_a$) to be 6,200 and 1,900 for protein-coding exons and CNEs, respectively. The

399    proportions of mutations with these selection coefficients were $9 \times 10^{-6}$ and $3 \times 10^{-4}$, respectively.

400    The more mildly beneficial class of mutations inferred for protein-coding exons and CNEs had

401    scaled effects of 210 and 7.0, respectively, and the proportion of mutations with these effects

402    were $3.5 \times 10^{-4}$ and $1.8 \times 10^{-2}$, respectively. In the case of CNEs, although two classes of

403    advantageous mutational effects gave the best fit to the data, the coefficient of variation for the

404    parameter estimates of the mildly selected class was large, and evidence for mildly beneficial

405    mutations is fairly weak in this case (Table 1).

406

407    The choice of recombination map strongly affected the estimated selection parameters obtained.

408    Use of the pedigree-based Cox map resulted in estimated selection coefficients that were typically

409    smaller than those obtained when assuming the LD-based recombination map (Supplementary

410    Table 3). This is because we found the troughs in diversity around both exons and CNEs were

411    shallower when calculating genetic distances using the pedigree-based map than when using the

412    LD-based map (Figure S1, S2).

413

414    BGS appears to contribute to the troughs in diversity around both protein-coding exons and CNEs

415    and causes an overall reduction in neutral diversity (Figure 4). Ignoring the contribution of BGS

416    (i.e. by setting *B* to 1.0 when fitting Equation 4 to the diversity troughs) resulted in a much poorer

417    model fit (Supplementary File 3). In the absence of BGS, the selection coefficients for

418    advantageous mutations required to explain the observed data are, as expected, far higher

419    (Supplementary File 3).

420

421

422

423

424  **Table 1** Parameters of positive selection in *M. m. castaneus* estimated by fitting a model of

425  selective sweeps and background selection to troughs in diversity around functional elements. The

426  frequency ($p_a$) and scaled selection coefficients ($\gamma_a$) for the two classes of advantageous effects are

427  given. Standard errors are shown in square brackets below point estimates.

428

429

| Element | $\gamma_{a,1}$ | $p_{a,1}$ | $\gamma_{a,2}$ | $p_{a,2}$ |
|---|---|---|---|---|
| Protein-Coding Exons | 6,170 | $0.80 \times 10^{-5}$ | 208 | $3.50 \times 10^{-4}$ |
| | [2,650] | [$0.50 \times 10^{-5}$] | [105] | [$1.80 \times 10^{-4}$] |
| CNEs | 1,910 | $1.30 \times 10^{-5}$ | 7.00 | $1.78 \times 10^{-2}$ |
| | [673] | [$0.60 \times 10^{-5}$] | [3.50] | [$1.29 \times 10^{-2}$] |

430

431  We did not include gene conversion events in our analysis, because gene conversion tracts, which

432  have an estimated mean length in mice of 135bp (Paigen et al. 2008), are relatively short

433  compared to the genetic distances we analysed (up to 100,000bp and 5,000bp for exons and CNEs,

434  respectively). Furthermore, the ratio of the rates of gene conversion and crossover events has

435  been estimated to be 0.105 in mice (Paigen et al. 2008). Overall, gene conversion is expected to

436  contribute little to the net frequency of recombination between neutral and selected sites.

437

438  The relative contribution of adaptive substitutions in protein-coding and regulatory regions

439  to fitness change in mice

440  An important goal of evolutionary biology is to understand the extent to which protein-coding and

441  regulatory elements contribute to phenotypic evolution (King and Wilson 1975; Wray 2007; Stern

442  and Orgogozo 2008; but see Hoekstra and Coyne 2007). Using our estimated selection parameters,

443  we can parameterise the following model of the rate of fitness change per generation ($\Delta W$)

444  brought about by the fixation of advantageous mutations. For a particular class of sites, assume

445  there are $\eta_a$ nucleotides in the genome at which new mutations occur at rate $\mu$ per nucleotide site

446  per generation. If the size of the breeding population is $N$, then $2N\mu$ new mutations enter the

447  population each generation. We assume that a proportion of the new mutations, $p_a$, is strongly

448  advantageous, with a selection coefficient of $s_a$ in heterozygous carriers. When the effectiveness

449  of selection exceeds that of genetic drift ($2N_e s_a > 1$), the fixation probability is approximately $2s_a$

450  (Haldane 1927). Once fixed, advantageous mutations increase population mean fitness by $s_a/h$,

451  where $h$ is the dominance coefficient, giving the following expression:

452

$$\Delta W = \frac{4N\mu p_a \eta_a s_a^2}{h} \tag{1}.$$

454

455 Since we are interested in the relative contribution to fitness change, and assumed that the

456 average point mutation rate is the same for CNEs and protein-coding exons, we can thus ignore $\mu$

457 in Equation 1. Note that the above model is conceptually similar to an approach taken by Lynch et

458 al. (1993) to model fitness change under mutational meltdown. We parametrized Equation 1 using

459 our estimated selection parameters. Note that we estimated two classes of beneficial mutational

460 effects for the two classes of functional elements. When parameterizing Equation 1, we summed

461 the fitness contributions over the two classes of fitness effect inferred for each element. We

462 calculated the ratio of $\Delta W$ for protein-coding exons and CNEs ($\Delta W_{Exons}/\Delta W_{CNEs}$) as a measure of

463 the relative contributions of the two types of elements to adaptive evolution (which also implicitly

464 assumes the same $h$ for all classes of mutation).

465

466 Our point estimates suggest that $\Delta W$ is larger for protein-coding regions than regulatory regions.

467 However, it is notable that the total genomic rate of fixation of beneficial mutations is higher for

468 CNEs than for coding regions (see also Halligan et al. 2013), but this reflects the fact that there are

469 approximately three times as many CNE bases as non-synonymous bases in the mouse genome.

470 Although the estimated genomic rate of fixation of beneficial mutations in CNEs is greater than

471 that of protein-coding exons (Table 2), the average strength of selection acting on a new

472 advantageous nonsynonymous mutation far exceeds that of CNEs (Table 2). Fitness change is

473 proportional to the square of the effect size, so that the change in population mean fitness

474 brought about by the fixation of advantageous mutations is substantially higher for protein-coding

475 exons than for CNEs. This result is sensitive to the choice of recombination map, since we inferred

476 stronger selection when assuming the LD-based map (Supplementary Table 3). Using a parametric

477 bootstrap approach, we found that $\Delta W_{Exons}/\Delta W_{CNEs}$ was significantly greater than 1 when using

478 the LD-based map, but not when assuming the pedigree-based map of Cox et al. (2009).

479

480 **Table 2** Estimates of the change in fitness brought about by the fixation of advantageous

481 mutations. Estimates were obtained assuming an effective populations size for *M. m. castaneus* of

482 420,000 and the selection parameters shown in Table 1.

483

| Recombination Map | Element | $\Delta W_{Exons}/\Delta W_{CNEs}$ | 95% Bootstrap interval |
|---|---|---|---|
| LD-based (*castaneus* map) | Protein-Coding Exons | 23.11 | 11.08 - 54.67 |
| | CNEs | | |
| Pedigree-based (Cox et al. 2009) | Protein-Coding Exons | 2.94 | 0.211 - 46.36 |
| | CNEs | | |

484

485

486    Selective sweeps and background selection in the mouse genome

487    The profiles of nucleotide diversity indicate the existence of pervasive effects of selection on

488    diversity across the genome (Figure 1, Figure 2). By fitting a model of sweeps to the troughs in

489    diversity around protein-coding exons and CNEs, while assuming that the troughs are partly

490    caused by BGS, we estimated the parameters of positively selected mutations occurring in the two

491    classes of element. Our analysis suggests that regulatory sequences experience a higher genomic

492    rate of newly arising advantageous mutations than protein-coding sites. However, the trough in

493    diversity around exons is both deeper and wider than what is observed around CNEs, and,

494    accordingly, we found that protein-coding regions experience more strongly selected mutations

495    than regulatory sequences. Using a different approach, Campos et al. (2017) came to a similar

496    conclusion for *D. melanogaster* by comparing UTRs with the coding sequences of genes.

497

498    Due to non-independence among the various *M. musculus* sub-species, we only estimated

499    parameters of positive selection for *M. m. castaneus,* the sub-species with the highest levels of

500    diversity.  Our selection parameter estimates for *M. m. castaneus* are fairly similar to estimates

501    obtained for European *M. m. domesticus* in an earlier study (Teschke et al. 2008).

502

503

### Limitations and next steps

505 There are a number of caveats concerning our estimates of positive selection parameters. Firstly,

506 we found that an exponential distribution of beneficial mutational effects provided a poorer fit to

507 the troughs in diversity compared to a model with two discrete classes of effects (Supplementary

508 File 3). However, the true DFE for advantageous mutations is almost certainly more complex than

509 the simple models assumed. The approach used in this study was based on average nucleotide

510 diversity across many sites, and we presumably had little power to infer a more complex model of

511 the DFE for advantageous mutations. Secondly, we have assumed that all elements of a particular

512 class share a common set of selection parameters. This is problematic, since CNEs could be

513 composed of several categories, such as promoters and enhancers, which may be subject to

514 different selective pressures. Indeed, different categories of protein-coding genes may also be

515 subject to different selection pressures. For example, immunity genes in *D. melanogaster*, virus

516 interacting proteins in humans and highly expressed genes in *Capsella grandiflora* appear to have

517 higher rates of adaptive substitutions than the respective genome wide averages (Enard et al.

518 2014; Obbard et al. 2009; Williamson et al. 2014). Thirdly, for a single class of advantageous

519 mutational effects, under the assumption that there is no interference among sweeps, the

520 predicted reductions in diversity caused by selective sweeps can be modelled as a simple

521 hyperbolic function (Equation 4). However, if the rate of sweeps is sufficiently high, and the rate of

522 recombination is sufficiently low, selective interference can cause the rate of sweeps to be lower

523 than predicted by a given strength of selection (Campos and Charlesworth 2019). This implies that

524 the strength of positive selection would be overestimated by our methods. A bias in the opposite

525 direction, which is likely to be more important for genomic regions with normal levels of

526 recombination, is caused by deviations from one of the assumptions underlying Equation 4, i.e.

527 that there is full recovery of nucleotide diversity between selective sweeps (Campos and

528 Charlesworth 2019; Charlesworth 2020). This would lead to the effects of sweeps to be

529 underpredicted. Incorporating more sophisticated models of selective sweeps into the inference

530 framework is a logical next step.

531

532 The architecture of functional elements in the mammalian genome is such that a single exon or

533 CNE is rarely far away from another functional element. When estimating the effects of sweeps on

534 neutral diversity, we excluded all putatively functional sites from our analysis windows, but

535 multiple linked elements may affect observed diversity at a given locus. For this reason, we did not

536     estimate the strength of positive selection acting on UTRs, although sweeps in these elements are

537     also likely to contribute to heterogeneity in $\pi$ across the mouse genome, as has been found in *D.*

538     *melanogaster* (Campos et al. 2017). The model fitted to the troughs in diversity assumes that

539     selection is generated by a single, idealised exon or CNE. However, there is variation in the length

540     of exons and CNEs across the genome. An analysis that models genome-wide heterogeneity in

541     diversity while taking into account the locations of individual functional elements, similar to the

542     method developed by Elyashiv et al. (2016) for *D. melanogaster*, could be a more powerful

543     approach. Note that the approach of Elyashiv et al. (2016) might not be applicable in all situations,

544     because it conditions the effects of sweeps on the locations of recent substitutions. In mice and

545     humans, patterns of diversity around nonsynonymous substitutions are indistinguishable from the

546     patterns of diversity around synonymous substitutions (Hernandez et al. 2011; Halligan et al.

547     2013). Developing a chromosome-wide analysis that conditions the effects of sweeps on the

548     locations of genomic elements rather than substitutions may be a useful avenue for further

549     research.

550

551     The model of sweeps we assumed involves positive selection acting on *de novo* mutations – the

552     so-called `hard', or `classic' sweep model. Studies in humans and *Drosophila* have, however,

553     suggested that 'soft' sweeps are common (Garud et al. 2015; Garud and Petrov 2016; Schrider and

554     Kern 2016; but see Harris et al. 2018). Soft selective sweeps occur when advantageous alleles

555     present in multiple copies in the population spread to fixation, which can occur if selection acts on

556     standing genetic variation or if multiple copies of the selected allele arise independently

557     (Hermisson and Pennings 2017). Additionally, adaptation acting on quantitative traits subject to

558     stabilising selection may generate partial sweeps, because changes in allele frequencies at many

559     loci can rapidly alter mean phenotypes, without necessarily causing fixations (Pritchard et al. 2010;

560     Jain and Stephan 2017) . The profiles of the reductions in diversity around soft and partial sweeps

561     differ from those expected under hard sweeps, and if either of the alternative types of sweep

562     were common, the assumption of a hard sweep model could result in spurious parameter

563     estimates (Elyashiv et al. 2016). Finally, the trough in diversity around a selective sweep in a

564     structured population is expected to be shallower than in a panmictic population because of the

565     longer time taken to reach fixation (Barton 2000; Santiago and Caballero 2005). If beneficial alleles

566     are frequently introduced via migration, we may therefore underestimate the strength of

567     selection.

568

569    Finally, it is important to note that CNEs are generally expected to represent regulatory sequences

570    that are deeply conserved. It has been demonstrated that the evolution of regulatory elements is

571    more dynamic than that of coding sequences, with major gains of new regulatory elements having

572    occurred in vertebrate and mammalian evolution (Mikkelsen et al. 2007; Lowe et al. 2011). If more

573    recently acquired regulatory elements, which may be absent from the CNE dataset, experience

574    stronger or more frequent adaptive substitutions, it is possible that we have underestimated the

575    contribution of regulatory changes to adaptive evolution. For instance, a recent gain of a new

576    regulatory element might have been caused by relatively strong positive selection acting on the

577    element as a whole, resulting in a single sweep event. This would fall outside the inference

578    framework developed here.

579

580    It seems likely that adaptation does not fit any one particular mode, but rather different functional

581    elements will be subject to a mixture of different types of sweep that may vary depending on the

582    genomic region. For example, adaptation may more commonly act on standing variation in

583    regulatory regions simply because they harbour greater nucleotide diversity than nonsynonymous

584    sites (Figure 2).

585

# Conclusions

586

587

588    In this study, we have shown that multiple wild mouse taxa exhibit patterns of genetic diversity

589    and divergence that are consistent with the action of natural selection. Furthermore, we have

590    shown that strong positive selection can explain the dips in diversity around protein-coding exons

591    and CNEs in *M. m. castaneus.* Finally, even though the framework we have adopted here is

592    incapable of distinguishing different modes of positive selection such as adaptation, sexual

593    selection and various forms of competition, the estimated parameters of positive selection

594    suggest that mutations in protein-coding regions may contribute more to the rate of change in

595    fitness under positive selection than regulatory mutations.

596

# Materials and Methods

## Genomic data

We re-analysed previously published genome sequences for the 54 wild-caught *Mus musculus* individuals described in Harr et al. (2016) and the 10 *M. m. castaneus* individuals and the *M. famulus* individual originally described in Halligan et al. (2010, 2013). The mouse samples belonged to three species: *Mus spretus, Mus musculus* and *Mus famulus*. The *M. spretus* individuals (*n* = 8) were from Madrid, Spain. The *M. musculus* individuals are composed of samples from the sub-species *M. m. domesticus, M. m. musculus* and *M. m. castaneus*. Three populations of *M. m. domesticus* were sampled (Massif Central, France, *n* = 8; Cologne-Bonn, Germany, *n* = 8; Ahvaz, Iran, *n* = 8) and three populations of *M. m. musculus* were sampled (Afghanistan, *n* = 6; Studenec, Czech Republic, *n* = 8; Mazar-e-Sharif, Kazakhstan, *n* = 8). We also analysed 10 *M. m. castaneus* described by Halligan et al. (2010, 2013), sampled in Himachal Pradesh, India. The one *M. famulus* individual, originated in Southern India, though Halligan et al. (2013) obtained it from the Montpellier Wild Mice Repository.

Harr et al. (2016) published and made available the variant calls obtained from the *M. musculus* samples described above in the form of VCF files. However, Harr et al. (2016) did not include invariant sites in their VCFs; for our purposes we required this information, so we re-called variants from their processed BAM files, available at http://wwwuser.gwdg.de/~evolbio/evolgen/wildmouse/. The data had been processed according to the GATK version 3 best practices pipeline, up to the step prior to variant calling. Briefly, all sequencing reads had been mapped to the *mm10* genome using *bwa-mem* (Li 2013). Reads were then sorted, merged and PCR duplicates were marked using *picardtools* (https://broadinstitute.github.io/picard/). Base Quality Score Recalibration was then applied using the dbSNP resource for mice (https://www.ncbi.nlm.nih.gov/snp) to produce analysis-ready alignments in BAM format. We generated BAM files for the *M. m castaneus* data and the *M. famulus* mice using the same procedure using FASTQ files downloaded from the European Nucleotide Archive (accession number PRJEB2176). For each of the mice, we called variants separately using the HaplotypeCaller tool from GATK3.7 (McKenna et al. 2010), with the options "– *emitRefConfidence BP_RESOLUTION  –max-alternate-alleles* 2", and made population-specific

25

627  VCF files using the GATK tools *combineGVCFs* and *genotypeGVCFs*. We restricted all analyses to

628  autosomal sites.

629

630  Outgroup information and CpG sites

631  In this study we used *M. famulus*, *Mus pahari* and *Rattus norvegicus* as the outgroup species.  For

632  each of the mouse taxa described above and each outgroup, we created a synthetic mm10-length

633  reference genome by replacing mm10 alleles with the major allele of the variant call set. In

634  addition, we constructed a synthetic genome for *R. norvegicus* by replacing mm10 alleles with the

635  homologous positions in the rat genome using the UCSC reciprocal best alignments between rn6

636  and mm10 (available at:

637  ftp://hgdownload.cse.ucsc.edu/goldenPath/rn6/vsMm10/reciprocalBest/) using custom Python

638  scripts. For an additional outgroup, more closely related to *Mus musculus* than the rat, we

639  obtained the homologous alleles from *Mus pahari* at mm10 positions using the ENSEMBL pairwise

640  alignments between the *M. pahari* reference sequence (Thybert et al. 2018) and mm10 (available

641  at: ftp://ftp.ensembl.org/pub/release-90/maf/ensembl-compara/pairwise_alignments/).

642

643  CpG sites have higher rates of spontaneous mutation than non-CpG sites, and identifying and

644  excluding CpG-prone sites is a conservative way of reducing the impact of CpG hypermutability on

645  analysis of population genomic data (Gaffney and Keightley 2008). For each of the rodent taxa, we

646  used the synthetic mm10-length reference genomes to identify the locations of CpG-prone sites,

647  defined as those sites in our synthetic references that were preceded by a C or followed by a G in

648  the 5' or 3' direction, respectively. All analyses presented in this paper excluded CpG-prone sites.

649

650  Annotations and identifying conserved non-coding elements

651  We downloaded the list of mouse-rat orthologs from https://www.ensembl.org/biomart/ and

652  extracted the annotations for each from version 38.93 of ENSEMBL

653  (Mus_musculus.GRCm38.93.gtf.gz; Howe et al. 2021). For each of the orthologs, we identified the

654  positions of 0-fold degenerate nonsynonymous and 4-fold degenerate synonymous sites using the

655  synthetic genomes for each of the mouse taxa and the outgroups described above. The locations

656  of 5' and 3' untranslated regions (UTRs) were retained for downstream analyses. We also retained

657  a list of all exonic positions in the mouse genome, regardless of orthology, for the purposes of

658  filtering out functionally constrained sites in downstream analyses.

26

659

660    There is evidence that synonymous sites within exonic splice enhancers (ESEs) in humans are

661    subject to purifying selection, and that ignoring ESEs can bias analyses that rely on the assumption

662    that synonymous sites evolve neutrally (Savisaar and Hurst 2018). Savisaar and Hurst (2018)

663    identified putative ESEs by comparing human gene sequences against various lists of ESE motifs.

664    They found that synonymous sites in regions matching ESE motifs had lower nucleotide diversity

665    than those outside of putative ESEs. We identified the locations of potential ESEs in protein-coding

666    genes orthologous between mice and rat using the merged list of ESEs described in Savisaar and

667    Hurst (2018) (kindly provided by Rosina Savisaar). For each of the mouse-rat orthologs, we

668    extracted the gene sequence and performed a string search against the list of ESE motifs. We

669    recorded the genomic position of each region matching an ESE motif and used them to filter out

670    the affected coding sites in downstream analysis.

671

672    We identified conserved non-coding elements (CNEs) in murid rodents using a 40-way alignment

673    of placental mammals downloaded from UCSC

674    (http://hgdownload.cse.ucsc.edu/goldenPath/mm10/multiz60way/). To avoid ascertainment bias,

675    the mouse and rat genomes in the 40-way alignment were converted to the character "N" prior to

676    calling conserved elements, following Williamson et al. (2014). We ran *phastCons* with the

677    following arguments --expected-length=45 --target-coverage=0.3 --rho=0.31. To identify CNEs, we

678    masked all exonic regions from the resulting file of *phastCons* elements using the complete list of

679    annotations from the 38.93 database (see above). The scripts and full pipeline used to identify

680    CNEs are available at https://github.com/rorycraig337/mouse_mm10_conserved_elements.

681

682    For each CNE identified in this way, we obtained the location of their flanking sequences, which

683    we used as neutral comparators in downstream analysis. For each CNE, we recorded the locations

684    of two loci of equal length upstream and downstream of the focal element, offset by 500bp. We

685    merged overlapping CNE-flanks and masked out sites that overlapped any CNE or exonic sites.

686

687    We analysed the mouse genomes assuming the pedigree-based genetic map of *Mus musculus*

688    constructed by Cox et al. (2009). The Cox map was constructed using data from 3,546 meioses

689    observed in crosses of common laboratory strains. The markers genotyped by Cox et al. (2009)

690    were mapped to the mm9 reference genome, but in the present study we converted the mm9

691    coordinates to mm10 positions as follows. The Cox map was downloaded from the Jackson

692    Laboratory website (http://cgd.jax.org/mousemapconverter/). The SNP positions of the Cox map

693    were then extracted and converted to mm10 positions using the online UCSC LiftOver tool

694    (https://genome.ucsc.edu/cgi-bin/hgLiftOver). The physical distances between the mm10 SNP

695    positions were then converted to units of genetic distance using the Jackson Laboratory's

696    conversion tool (http://cgd.jax.org/mousemapconverter/). We also analysed the mouse genomes

697    using an LD-based recombination map inferred from the sample of *M. m. castaneus* individuals, as

698    described in the Appendix.

699

700    Mouse analysis – Patterns of nucleotide diversity around selected sites

701    For each of the *M. musculus* sub-species and *Mus spretus*, we examined patterns of nucleotide

702    diversity around protein-coding exons and CNEs. From the edges of protein-coding exons (CNEs),

703    polymorphism data and divergence from the rn6 rat reference genome were extracted in windows

704    of 1Kbp (100bp) extending to distances of 100Kbp (5Kbp). Analysis windows only extended to the

705    midway point between adjacent elements. Sites within the exons of protein-coding genes or CNEs

706    were excluded from analysis windows. The genetic distance between an analysis window and a

707    focal element was calculated either from the pedigree-based genetic map constructed using

708    common lab strains of *M. musculus* (Cox et al. 2009) or the linkage disequilibrium (LD) based

709    recombination map for *M. m. castaneus*. The SFS and divergence from were recorded for each

710    analysis window. Analysis windows were then binned based on the genetic distance from the focal

711    element, and the SFS and divergence from individual windows were collated. Because LD-based

712    and pedigree-based recombination maps have different features and shortcomings (see Results),

713    we performed analyses based on both genetic maps.

714

715    Estimating the unfolded site frequency spectrum, summary statistics and the

716    distribution of fitness effects

717    We analysed genetic variation for five different classes of sites in the genome, i.e. the 0-fold and 4-

718    fold degenerate sites and UTRs of protein-coding genes, CNEs and CNE-flanks.  For each class of

719    sites, we inferred the unfolded site frequency spectrum (uSFS), which is the distribution of derived

720    allele frequencies in the mouse samples. The uSFS was inferred by maximum likelihood using the

721    two-outgroup method of Keightley and Jackson (2018) using *Mus famulus* and *Mus pahari* as

722    outgroups. We compared the fit of 1-, 2- and 6- parameter mutation rate models using *est-sfs*

723    (v2.03; Keightley and Jackson 2018). Consistently, a model with 6 mutation rate parameters (i.e.

724    the R6 model from Keightley and Jackson 2018) provided the best fit to the data (as judged by

725    model likelihoods), but the uSFS and lineage specific divergences that were estimated under the 2-

726    parameter and 6-parameter models were almost identical in all cases (Supplementary File 1), so

727    we have used the results from the 2-parameter model in our analyses for all taxa. For each taxon

728    and class of sites, we performed 100 bootstraps, sampling genes or CNEs with replacement. We

729    inferred the uSFS for each bootstrapped dataset as above. For each class of sites, we calculated

730    nucleotide diversity ($\pi$) and Tajima's $D$ from the inferred uSFS for each bootstrap sample.

731

732    Estimates of the distribution of fitness effects (DFE) were obtained by analysing the unfolded site

733    frequency spectrum for each class of functional site using *polyDFE* v2 (Tataru and Bataillon 2019).

734    For 0-fold sites and UTRs, we used 4-fold degenerate synonymous sites as the neutral comparator,

735    and for CNEs we used CNE-flanks. Using *polyDFE2*, we fitted a gamma distribution of deleterious

736    mutations effects and a single class of beneficial mutations (using the -model B option). We

737    excluded between species divergence from the analysis using the "-w" option. We fitted the uSFS

738    data for each of the bootstrap replicates described above.

739

740    Simulating background selection

741    There is substantial evidence that background selection (BGS) contributes to troughs in diversity

742    around protein-coding exons and CNEs (Halligan et al. 2013; Booker and Keightley 2018). For our

743    analysis, we therefore required estimates of the effect of BGS on neutral diversity, $B$, at varying

744    distances from functional elements. Estimates of $B$ were included as covariates when fitting a

745    model of selection at linked sites. However, when purifying selection is weak ($\gamma_d < 5$) analytical

746    formulae for calculating $B$ over-predict the effects of BGS (Gordo et al. 2002; Good et al. 2014),

747    and weakly deleterious mutations appear to comprise a large fraction of the DFEs for mice (Figure

748    3 and Halligan et al. 2013). We therefore opted to obtain estimates of the variation in $B$ from

749    forward-in-time simulations that modelled the entire range of fitness effects inferred for mice.

750

751    We used *SLiM* v3.2 (Haller and Messer 2019) for this purpose. Following Booker and Keightley

752    (2018), we incorporated the actual distribution of functional elements (the coding exons and UTRs

753    of protein-coding genes, and CNEs) and the estimated recombination rates. 1 Mbp regions of the

754    mouse genome were randomly sampled and the functional annotations in the sampled regions

755    were used as the basis of a simulation replicate. The parameters of the gamma distributions of

756    fitness effects for deleterious mutations estimated for 0-fold sites, UTRs and CNEs were used in

757    the simulations for the respective elements. The recombination rate variation present in the

758    sampled region of the mouse genome was included in the simulations using either the pedigree-

759    based map from Cox et al. (2009) or the LD-based recombination map for *M. m. castaneus*. When

760    assuming the Cox map, the recombination rates (in units of cM/Mbp) were scaled in the

761    simulations by a factor of $420r$, assuming $N_e = 420,000$ for wild *M. m. castaneus*. In the case of the

762    LD-based estimates of the recombination rate, population-scaled recombination rates (in units of

763    $4N_er$) were simply divided by $4N$, where $N$ was the simulated population size. Populations of $N =$

764    1,000 diploid individuals were simulated for 20,000 generations. We set the mutation rate such

765    that the neutral expectation $\pi_0 = 4N_e\mu = 0.01$, based on the upper estimate of nucleotide diversity

766    observed in the *M. m. castaneus* genome (Figure 1). Given the simulated population size of 1,000

767    diploids, $4N_e\mu = 0.01$ corresponded to a point mutation rate of $\mu = 2.5 \times 10^{-6}$. We used the tree-

768    sequence recording option in *SLiM* to record the genealogies of the simulated populations, so

769    modelling neutral mutations in *SLiM* was not required. Instead, neutral mutations were added to

770    the recorded coalescent trees at a rate $\mu$ using *PySLiM*

771    (https://pyslim.readthedocs.io/en/latest/introduction.html). We sampled 200 haploid

772    chromosomes from the population and extracted $B = \pi/\pi_0$ as a function of genetic distance from

773    both protein-coding exons and CNEs. Data were extracted from the simulated populations in the

774    same way as for the empirical data. To obtain smoothed $B$ values, we fitted LOESS curves to the

775    average $\pi$ observed around functional elements in the simulated data. We fitted LOESS curves

776    using a span parameter of 0.3 and the number of sites contributing to each analysis bin as weights

777    in R (v3.4.2).

778

779    Model of recurrent selective sweeps and background selection

780    Our analysis is a modification of that of  Elyashiv et al. (2016) and Campos et al. (2017), where

781    expressions were described for the neutral diversity expected under the combined effects of BGS

782    and sweeps. Consider a haplotype with a set of neutral sites linked to a site that is the target of

783    positive selection. A new, semi-dominant advantageous mutation with heterozygous selection

784    coefficient $s_a$ occurs at site $i$ on the haplotype and spreads to fixation. Recombination between

785    sites $i$ and $k$ uncouples the neutral and selected sites at rate $r_{i,k}$ per generation.  The expected

786    change in neutral diversity at site $k$ ($\Delta\pi_k$), relative to its expectation in the absence of selection ($\pi_0$)

787    is given by

788    $$\frac{\Delta\pi_k}{\pi_0} = -(4N_es_a)^{\frac{-2r_{i,k}}{s_a}}.$$    (2)

30

789

790    See Barton (2000), Charlesworth and Charlesworth (2010, p411) or Campos and Charlesworth

791    (2019) for derivations of Equation 2. This approximation assumes that selection pressure on the

792    advantageous allele satisfies $N_e s_a \gg 1$, so that the sweep can be treated deterministically

793    following an initial stochastic establishment phase. Under this assumption, the quantity $-\Delta \pi_k$ in

794    Equation 1 can be equated to the probability of a sweep-induced coalescent event at site $k$ (Wiehe

795    and Stephan 1993). For a particular class of functional elements (e.g. protein-coding exons),

796    sweeps occur at a rate of $V_a = 2\mu p_a \gamma_a$ per nucleotide per generation (Kimura and Ohta 1971),

797    where $\mu$ is the mutation rate per nucleotide site, $p_a$ is the proportion of new mutations that are

798    advantageous and $\gamma_a$ is the scaled selection coefficient ($4N_e s_a$) for these mutations. If $V_a$ is

799    sufficiently low, such that sweeps do not interfere with each other, the total probability of sweep-

800    induced coalescence for a neutral site caused by selection at a linked functional element is:

801

802    $$P_{sc,k} = V_a \tau \gamma_a^{\frac{-2r_{i,k}}{s_a}} , \qquad\qquad (3)$$

803

804    where $\tau$ is the number of sites in a particular class of functional element. In our analysis of data

805    from wild mice, $r_{i,k}$ was measured from the end of a functional element to the centre of an analysis

806    window.

807

808    The effects of background selection at site $k$ can be represented by multiplying the effective

809    population size by a factor $B_k$. The probability of coalescence for a neutral allele affected by BGS is

810    thus $1/(2B_k N_e)$. We assume that coalescent events caused by BGS and sweeps follow independent

811    exponential distributions, so thar the rate of coalescence induced by the two processes is the sum

812    of $1/(2B_k N_e)$ and $P_{sc,k}$. We also multiply the sweep effect $P_{sc,k}$ by $B_k$ to reflect the reduction in the

813    fixation probability of a new advantageous mutation as a result of the reduction in $N_e$ caused by

814    BGS, following Kim and Stephan (2000). Simulations show that this may overestimate the effect of

815    BGS on fixation probabilities (Campos and Charlesworth 2019), we thus compared selection

816    parameters with and without including background selection.

817

818    Writing the reciprocal of the rate of coalescence at the neutral site under the combined effects of

819    BGS and sweeps as $T_k$ (which is equivalent to the expected time to coalescence of a pair of alleles),

820    and expressing it relative to the expected time to coalescence under neutrality ($T_0$), we have:

821

$$\frac{T_k}{T_0} \approx \frac{\pi_k}{\pi_0} \approx \frac{1}{B_k^{-1} + 2N_e B_k P_{sc,k}}. \tag{4}$$

823

824 We estimated parameters of advantageous mutations by fitting Equation 4 to the relationship

825 between nucleotide diversity and genetic distance from functional elements, using non-linear

826 least squares with the *lmfit* (v0.9.7) package for Python 2.7. When modelling a single class of

827 fitness effects, we estimated $\gamma_a$ and $p_a$ using Equation 4. To incorporate two discrete classes of

828 advantageous mutational effects, we modified Equation 4, replacing $P_{sc,k}$ in Equation 4 with

829

$$P_{sc,k} = V_{a,2}\tau\gamma_{a,2}^{\frac{-2r_{i,k}}{s_{a,2}}} + V_{a,2}\tau\gamma_{a,2}^{\frac{-2r_{i,k}}{s_{a,2}}}, \tag{5}$$

831

832 where the subscripts 1 and 2 refer to the two different classes of fitness effects.

833

834 To model an exponential distribution of fitness effects, we replaced $P_{sc,k}$ with

835

$$P_{sc,k} = \int_0^\infty V_a\tau\gamma_a^{\frac{-2r_{i,k}}{s_a}} \phi(\gamma_a|\bar{\gamma}_a)\mathrm{d}\gamma_a, \tag{6}$$

837

838 where $\phi(\gamma_a|\bar{\gamma}_a)$ is the probability density function of an exponential distribution with mean $\bar{\gamma}_a$. In

839 all cases, we used the average length of protein-coding exons or CNEs, 152.0 and 50.0

840 respectively, as $\tau$ when fitting equation 4. We assumed that $N_e$ = 426,200, based on $4N_e\mu$ = 0.0092

841 and a mutation rate of 5.4 x $10^{-9}$ (Uchimura et al. 2015).

842

843 We used our estimates of positive selection parameters to quantify the relative contributions of

844 positive selection in protein-coding exons versus CNEs to the change in population mean fitness.

845 To obtain confidence intervals around our estimates of the ratio of fitness contributed by positive

846 selection in exons versus CNEs ($\Delta W_{Exons}/\Delta W_{CNEs}$), we used a parametric bootstrap approach. For

847 each estimated $\gamma_a$ and $p_a$ parameter, we sampled random values from a truncated normal

848 distribution with mean equal to the parameter estimate and variance equal to the square of the

849 standard error of the parameter estimate. The truncated normal distribution had a lower bound of

850 0.0, since values of $\gamma_a$ and $p_a$ below 0 are biologically impossible. For the calculation of $\Delta W_{Exons}/$

851 $\Delta W_{CNEs}$, we performed 1,000 bootstraps and used them to estimate 95% confidence intervals.

32

852

856

# Author Contributions

858    TRB, BC and PDK devised the study. TRB, BCJ, RJC analysed the data. TRB wrote the first draft of

859    the manuscript. All authors contributed to the writing and editing of the manuscript.

860

# Acknowledgements

868

# References

Barton, N. H. (2000). Genetic hitchhiking. *Philos Trans R Soc Lond B Biol Sci*, *355*(1403), 1553–1562. https://doi.org/10.1098/rstb.2000.0716

Booker, T R, Ness, R. W., & Keightley, P. D. (2017). The recombination landscape in wild house mice inferred using population genomic data. *Genetics*, *207*(1), 297–309. https://doi.org/10.1534/genetics.117.300063

Booker, Tom R. (2020). Inferring Parameters of the Distribution of Fitness Effects of New Mutations When Beneficial Mutations Are Strongly Advantageous and Rare. *G3*, *10*(7), 2317–2326. https://doi.org/10.1534/g3.120.401052

Booker, Tom R, & Keightley, P. D. (2018). Understanding the Factors That Shape Patterns of Nucleotide Diversity in the House Mouse Genome. *Molecular Biology and Evolution*, *35*(12), 2971–2988.

Campos, J L, Zhao, L., & Charlesworth, B. (2017). Estimating the parameters of background selection and selective sweeps in Drosophila in the presence of gene conversion. *Proceedings of the National Academy of Sciences*, *114*(24), E4762–E4771. https://doi.org/10.1073/pnas.1619434114 10.5061/dryad.vs264)

Campos, José Luis, & Charlesworth, B. (2019). The effects on neutral variability of recurrent selective sweeps and background selection. *Genetics*, *212*(1), 287–303. https://doi.org/10.1534/genetics.119.301951

Carroll, S. B. (2005). Evolution at two levels: on genes and form. *PLoS Biology*, *3*(7), e245. https://doi.org/10.1371/journal.pbio.0030245

Chan, A. H., Jenkins, P. A., & Song, Y. S. (2012). Genome-wide fine-scale recombination rate variation in Drosophila melanogaster. *PLoS Genetics*, *8*(12), e1003090. https://doi.org/10.1371/journal.pgen.1003090

Charlesworth, B, & Charlesworth, D. (2010). *Elements of Evolutionary Genetics*. Roberts \& Company.

Charlesworth, Brian. (2012). The effects of deleterious mutations on evolution at linked sites. *Genetics*, *190*(1), 5–22. https://doi.org/10.1534/genetics.111.134288

Charlesworth, Brian. (2020). How good are predictions of the effects of selective sweeps on levels of neutral diversity? *Genetics*, *216*(4), 1217–1238. https://doi.org/10.1534/genetics.120.303734

Clark, A. G., Wang, X., & Matise, T. (2010). Contrasting Methods of Quantifying Fine Structure of

901    Human Recombination. *Annual Review of Genomics and Human Genetics*, *11*(1), 45–64.

902    https://doi.org/10.1146/annurev-genom-082908-150031

903  Comeron, J. (2014). Background selection as a baseline for nucleotide variation across the

904    Drosophila genome. *PLoS Genetics*, *10*(6). https://doi.org/10.1371/

905  Cox, A., Ackert-Bicknell, C. L., Dumont, B. L., Ding, Y., Bell, J. T., Brockmann, G. A., Wergedal, J. E.,

906    Bult, C., Paigen, B., Flint, J., Tsaih, S. W., Churchill, G. A., & Broman, K. W. (2009). A new

907    standard genetic map for the laboratory mouse. *Genetics*, *182*(4), 1335–1344.

908    https://doi.org/10.1534/genetics.109.105486

909  Delaneau, O., Howie, B., Cox, A. J., Zagury, J. F., & Marchini, J. (2013). Haplotype estimation using

910    sequencing reads. *Am J Hum Genet*, *93*(4), 687–696.

911    https://doi.org/10.1016/j.ajhg.2013.09.002

912  Dumont, B. L., & Payseur, B. A. (2011). Genetic analysis of genomic-scale recombiantion rate

913    evolution in house mice. *PLoS Genetics*, *7*(6), 11. https://doi.org/10.1371/

914  Elyashiv, E., Sattath, S., Hu, T. T., Strutsovsky, A., McVicker, G., Andolfatto, P., Coop, G., & Sella, G.

915    (2016). A genomic map of the effects of linked selection in Drosophila. *PLoS Genetics*, *12*(8),

916    e1006130. https://doi.org/10.1371/journal.pgen.1006130

917  Enard, D., Messer, P. W., & Petrov, D. A. (2014). Genome-wide signals of positive selection in

918    human evolution. *Genome Research*, *24*(6), 885–895. https://doi.org/10.1101/gr.164822.113

919  Ewing, G. B., & Jensen, J. D. (2016). The consequences of not accounting for background selection

920    in demographic inference. *Molecular Ecology*, *25*(1), 135–141.

921    https://doi.org/10.1111/mec.13390

922  Falconer, D. S., & Mackay, T. F. C. (1996). *Introduction to Quantitative Genetics (Fourth Edition)*

923    (Fourth). Pearson Education Limited.

924  Fisher, R. A. (1918). The Correlation Between Relatives on the Supposition of Mendelian

925    Inheritance. *Transactions of the Royal Society of Edinburgh*, *52*, 399–433.

926  Gaffney, D. J., & Keightley, P. D. (2008). Effect of the assignment of ancestral CpG state on the

927    estimation of nucleotide substitution rates in mammals. *BMC Evol Biol*, *8*, 265.

928    https://doi.org/10.1186/1471-2148-8-265

929  Garud, N. R., Messer, P. W., Buzbas, E. O., & Petrov, D. A. (2015). Recent selective sweeps in North

930    American Drosophila melanogaster show signatures of soft sweeps. *PLoS Genetics*, *11*(2),

931    e1005004. https://doi.org/10.1371/journal.pgen.1005004

932  Garud, N. R., & Petrov, D. A. (2016). Elevated linkage disequilibrium and signatures of soft sweeps

933    are common in Drosophila melanogaster. *Genetics*, *203*(2), 863–880.

35

934      https://doi.org/10.1534/genetics.115.184002

935  Geraldes, A., Basset, P., Smith, K. L., & Nachman, M. W. (2011). Higher differentiation among

936      subspecies of the house mouse (Mus musculus) in genomic regions with low recombination.

937      *Molecular Ecology*, *20*(22), 4722–4736. https://doi.org/10.1111/j.1365-294X.2011.05285.x

938  Good, B. H., Walczak, A. M., Neher, R. A., & Desai, M. M. (2014). Genetic Diversity in the

939      Interference Selection Limit. *PLoS Genetics*, *10*(3).

940      https://doi.org/10.1371/journal.pgen.1004222

941  Gordo, I., Navarro, A., & Charlesworth, B. (2002). Muller's Ratchet and the Pattern of Variation at a

942      Neutral Locus. *Genetics*, *161*, 835–848.

943  Haldane, J. B. S. (1927). A Mathematical Theory of Natural and Artificial Selection, Part V: Selection

944      and Mutation. *Mathematical Proceedings of the Cambridge Philosophical Society*, *23*(7), 838–

945      844. https://doi.org/10.1017/S0305004100015644

946  Haller, B. C., & Messer, P. W. (2019). SLiM 3: Forward Genetic Simulations Beyond the Wright-

947      Fisher Model. *Molecular Biology and Evolution*, *36*(3), 632–637.

948  Halligan, D. L., Kousathanas, A., Ness, R. W., Harr, B., Eory, L., Keane, T. M., Adams, D. J., &

949      Keightley, P. D. (2013). Contributions of protein-coding and regulatory change to adaptive

950      molecular evolution in murid rodents. *PLoS Genetics*, *9*(12), e1003995.

951      https://doi.org/10.1371/journal.pgen.1003995

952  Halligan, D. L., Oliver, F., Eyre-Walker, A., Harr, B., & Keightley, P. D. (2010). Evidence for pervasive

953      adaptive protein evolution in wild mice. *PLoS Genetics*, *6*(1), e1000825.

954      https://doi.org/10.1371/journal.pgen.1000825

955  Harr, B., Karakoc, E., Neme, R., Teschke, M., Pfeifle, C., Pezer, Ž., Babiker, H., Linnenbrink, M.,

956      Montero, I., Scavetta, R., Abai, M. R., Molins, M. P., Schlegel, M., Ulrich, R. G., Altmüller, J.,

957      Franitza, M., Büntge, A., Künzel, S., & Tautz, D. (2016). Genomic resources for wild

958      populations of the house mouse, Mus musculus and its close relative Mus spretus. *Scientific*

959      *Data*, *3*. https://doi.org/10.1038/sdata.2016.75

960  Harris, R. B., Sackman, A., & Jensen, J. D. (2018). On the unfounded enthusiasm for soft selective

961      sweeps II: Examining recent evidence from humans, flies, and viruses. *PLoS Genetics*, *14*(12),

962      e1007859. https://doi.org/10.1371/journal.pgen.1007859

963  Hermisson, J., & Pennings, P. S. (2017). Soft sweeps and beyond: understanding the patterns and

964      probabilities of selection footprints under rapid adaptation. *Methods in Ecology and*

965      *Evolution*, *8*(6), 700–716. https://doi.org/10.1111/2041-210x.12808

966  Hernandez, R. D., Kelly, J. L., Elyashiv, E., Melton, S. C., Auton, A., McVean, G., Project, 1000

967  Genomes, Sella, G., & Przeworski, M. (2011). Classic selective sweeps were rare in recent

968  human evolution. *Science*, *331*, 920–924.

969  Hoekstra, H. E., & Coyne, J. A. (2007). The locus of evolution: Evo devo and the genetics of

970  adaptation. *Evolution*, *61*(5), 995–1016. https://doi.org/10.1111/j.1558-5646.2007.00105.x

971  Howe, K. L., Achuthan, P., Allen, J., Allen, J., Alvarez-Jarreta, J., Ridwan Amode, M., Armean, I. M.,

972  Azov, A. G., Bennett, R., Bhai, J., Billis, K., Boddu, S., Charkhchi, M., Cummins, C., da Rin

973  Fioretto, L., Davidson, C., Dodiya, K., El Houdaigui, B., Fatima, R., … Flicek, P. (2021). Ensembl

974  2021. *Nucleic Acids Research*, *49*(D1), D884–D891. https://doi.org/10.1093/nar/gkaa942

975  Jain, K., & Stephan, W. (2017). Rapid adaptation of a polygenic trait after a sudden environmental

976  shift. *Genetics*, *206*(1), 389–406.

977  Johri, P., Riall, K., Becher, H., Excoffier, L., Charlesworth, B., & Jensen, J. D. (2021). The Impact of

978  Purifying and Background Selection on the Inference of Population History: Problems and

979  Prospects. *Molecular Biology and Evolution*. https://doi.org/10.1093/molbev/msab050

980  Keightley, P. D., & Jackson, B. C. (2018). Inferring the Probability of the Derived vs. the Ancestral

981  Allelic State at a Polymorphic Site. *Genetics*, *209*(3), 897–906.

982  Kern, A. D., & Hahn, M. W. (2018). The neutral theory in light of natural selection. *Molecular*

983  *Biology and Evolution*, *35*(6), 1366–1371. https://doi.org/10.1093/molbev/msy092

984  Kim, Y., & Stephan, W. (2000). Joint effects of genetic hitchhiking and background selection on

985  neutral variation. *Genetics*, *155*, 1415–1427.

986  Kimura, M., & Ohta, T. (1971). *Theoretical aspects of population genetics*. Princeton Univ. Press.

987  King, M.-C., & Wilson, A. C. (1975). Evolution at two levels in humans and chimpanzees. *Science*,

988  *188*(4184), 107–116.

989  Kousathanas, A., & Keightley, P. D. (2013). A comparison of models to infer the distribution of

990  fitness effects of new mutations. *Genetics*, *193*(4), 1197–1208.

991  https://doi.org/10.1534/genetics.112.148023

992  Lawal, R. A., Arora, U. P., & Dumont, B. L. (2021). Selection shapes the landscape of functional

993  variation in wild house mice 1 2. *BioRxiv*, 2021.05.12.443838.

994  https://doi.org/10.1101/2021.05.12.443838

995  Leffler, E. M., Bullaughey, K., Matute, D. R., Meyer, W. K., Segurel, L., Venkat, A., Andolfatto, P., &

996  Przeworski, M. (2012). Revisiting an old riddle: what determines genetic diversity levels

997  within species? *PLoS Biology*, *10*(9), e1001388. https://doi.org/10.1371/journal.pbio.1001388

998  Li, H. (2013). *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM*.

999  http://arxiv.org/abs/1303.3997

37

1000  Lindblad-Toh, K., Garber, M., Zuk, O., Lin, M. F., Parker, B. J., Washietl, S., Kheradpour, P., Ernst, J.,
1001      Jordan, G., Mauceli, E., Ward, L. D., Lowe, C. B., Holloway, A. K., Clamp, M., Gnerre, S., Alföldi,
1002      J., Beal, K., Chang, J., Clawson, H., … Sodergren, E. (2011). A high-resolution map of human
1003      evolutionary constraint using 29 mammals. *Nature*, *478*(7370), 476–482.
1004      https://doi.org/10.1038/nature10530

1005  Lowe, C. B., Kellis, M., Siepel, A., Raney, B. J., Clamp, M., Salama, S. R., Kingsley, D. M., Lindblad-
1006      Toh, K., & Haussler, D. (2011). Three periods of regulatory innovation during vertebrate
1007      evolution. *Science*, *333*(6045), 1019–1024. https://doi.org/10.1126/science.1202702

1008  Lynch, M., Bürger, R., Butcher, D., & Gabriel, W. (1993). The Mutational Meltdown in Asexual
1009      Populations. *Journal of Heredity*, *84*(5), 339–344.
1010      https://doi.org/10.1093/oxfordjournals.jhered.a111354

1011  McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K.,
1012      Altshuler, D., Gabriel, S., Daly, M., & DePristo, M. A. (2010). The genome analysis toolkit: A
1013      MapReduce framework for analyzing next-generation DNA sequencing data. *Genome*
1014      *Research*, *20*(9), 1297–1303. https://doi.org/10.1101/gr.107524.110

1015  McVicker, G., Gordon, D., Davis, C., & Green, P. (2009). Widespread genomic signatures of natural
1016      selection in hominid evolution. *PLoS Genetics*, *5*(5), e1000471.
1017      https://doi.org/10.1371/journal.pgen.1000471

1018  Mikkelsen, T. S., Wakefield, M. J., Aken, B., Amemiya, C. T., Chang, J. L., Duke, S., Garber, M.,
1019      Gentles, A. J., Goodstadt, L., Heger, A., Jurka, J., Kamal, M., Mauceli, E., Searle, S. M. J.,
1020      Sharpe, T., Baker, M. L., Batzer, M. A., Benos, P. V., Belov, K., … Lindblad-Toh, K. (2007).
1021      Genome of the marsupial Monodelphis domestica reveals innovation in non-coding
1022      sequences. *Nature*, *447*(7141), 167–177. https://doi.org/10.1038/nature05805

1023  Obbard, D. J., Welch, J. J., Kim, K. W., & Jiggins, F. M. (2009). Quantifying adaptive evolution in the
1024      Drosophila immune system. *PLoS Genetics*, *5*(10), 1000698.
1025      https://doi.org/10.1371/journal.pgen.1000698

1026  Paigen, K., Szatkiewicz, J. P., Sawyer, K., Leahy, N., Parvanov, E. D., Ng, S. H., Graber, J. H., Broman,
1027      K. W., & Petkov, P. M. (2008). The recombinational anatomy of a mouse chromosome. *PLoS*
1028      *Genetics*, *4*(7), e1000119. https://doi.org/10.1371/journal.pgen.1000119

1029  Pouyet, F., Aeschbacher, S., Thiéry, A., & Excoffier, L. (2018). Background selection and biased
1030      gene conversion affect more than 95% of the human genome and bias demographic
1031      inferences. *ELife*, *7*. https://doi.org/10.7554/eLife.36317

1032  Pritchard, J. K., Pickrell, J. K., & Coop, G. (2010). The genetics of human adaptation: hard sweeps,

1033    soft sweeps, and polygenic adaptation. *Current Biology*, *20*(4), R208-15.

1034    https://doi.org/10.1016/j.cub.2009.11.055

1035    Santiago, E., & Caballero, A. (2005). Variation after a selective sweep in a subdivided population.

1036    *Genetics*, *169*(1), 475–483. https://doi.org/10.1534/genetics.104.032813

1037    Sattath, S., Elyashiv, E., Kolodny, O., Rinott, Y., & Sella, G. (2011). Pervasive adaptive protein

1038    evolution apparent in diversity patterns around amino acid substitutions in Drosophila

1039    simulans. *PLoS Genetics*, *7*(2), e1001302. https://doi.org/10.1371/journal.pgen.1001302

1040    Savisaar, R., & Hurst, L. D. (2018). Exonic splice regulation imposes strong selection at synonymous

1041    sites. *Genome Research*, *28*(10), 1442–1454. https://doi.org/10.1101/gr.233999.117

1042    Schrider, D. R., & Kern, A. D. (2016). S/HIC: Robust Identification of Soft and Hard Sweeps Using

1043    Machine Learning. *PLoS Genetics*, *12*(3), e1005928.

1044    https://doi.org/10.1371/journal.pgen.1005928

1045    Siepel, A., Bejerano, G., Pedersen, J. S., Hinrichs, A. S., Hou, M., Rosenbloom, K., Clawson, H.,

1046    Spieth, J., Hillier, L. D. W., Richards, S., Weinstock, G. M., Wilson, R. K., Gibbs, R. A., Kent, W.

1047    J., Miller, W., & Haussler, D. (2005). Evolutionarily conserved elements in vertebrate, insect,

1048    worm, and yeast genomes. *Genome Research*, *15*(8), 1034–1050.

1049    https://doi.org/10.1101/gr.3715005

1050    Simonsen, K. L., Churchill, G. A., & Aquadro, C. F. (1995). Properties of statistical tests of neutrality

1051    for DNA polymorphism data. *Genetics*, *141*(1).

1052    Stephan, W. (2019). Selective sweeps. *Genetics*, *211*(1), 5–13.

1053    https://doi.org/10.1534/genetics.118.301319

1054    Stern, D. L., & Orgogozo, V. (2008). The loci of evolution: How predictable is genetic evolution?

1055    *Evolution*, *62*(9), 2155–2177. https://doi.org/10.1111/j.1558-5646.2008.00450.x

1056    Tataru, P., & Bataillon, T. (2019). polyDFEv2.0: testing for invariance of the distribution of fitness

1057    effects within and across species. *Bioinformatics*, *35*(16), 2868–2869.

1058    https://doi.org/10.1093/bioinformatics/bty1060

1059    Tataru, P., Mollion, M., Glemin, S., & Bataillon, T. (2017). Inference of distribution of fitness effects

1060    and proportion of adaptive substitutions from polymorphism data. *Genetics*, *207*(3), 1103–

1061    1119. https://doi.org/10.1534/genetics.117.300323

1062    Teschke, M., Mukabayire, O., Wiehe, T., & Tautz, D. (2008). Identification of selective sweeps in

1063    closely related populations of the house mouse based on microsatellite scans. *Genetics*, *180*,

1064    1537–1545.

1065    Thybert, D., Roller, M., Navarro, F. C. P., Fiddes, I., Streeter, I., Feig, C., Martin-Galvez, D.,

1066      Kolmogorov, M., Janoušek, V., Akanni, W., Aken, B., Aldridge, S., Chakrapani, V., Chow, W.,

1067      Clarke, L., Cummins, C., Doran, A., Dunn, M., Goodstadt, L., … Flicek, P. (2018). Repeat

1068      associated mechanisms of genome evolution and function revealed by the Mus caroli and

1069      Mus pahari genomes. *Genome Research*, *28*(4), 448–459.

1070      https://doi.org/10.1101/gr.234096.117

1071      Torres, R., Stetter, M. G., Hernandez, R. D., & Ross-Ibarra, J. (2020). The Temporal Dynamics of

1072      Background Selection in Nonequilibrium Populations. *Genetics*, *214*(4), 1019–1030.

1073      https://doi.org/10.1534/genetics.119.302892

1074      Uchimura, A., Higuchi, M., Minakuchi, Y., Ohno, M., Toyoda, A., Fujiyama, A., Miura, I., Wakana, S.,

1075      Nishino, J., & Yagi, T. (2015). Germline mutation rates and the long-term phenotypic effects

1076      of mutation accumulation in wild-type laboratory mice and mutator mice. *Genome Research*,

1077      *25*(8), 1125–1134. https://doi.org/10.1101/gr.186148.114

1078      Wiehe, T., & Stephan, W. (1993). Analysis of a genetic hitchhiking model, and its application to

1079      DNA polymorphism data from Drosophila melanogaster. *Molecular Biology and Evolution*,

1080      *10*(4), 842–854.

1081      Williamson, R. J., Josephs, E. B., Platts, A. E., Hazzouri, K. M., Haudry, A., Blanchette, M., & Wright,

1082      S. I. (2014). Evidence for widespread positive and negative selection in coding and conserved

1083      noncoding regions of Capsella grandiflora. *PLoS Genetics*, *10*(9), e1004622.

1084      https://doi.org/10.1371/journal.pgen.1004622

1085      Wray, G. A. (2007). The evolutionary significance of cis-regulatory mutations. *Nature Reviews*

1086      *Genetics*, *8*(3), 206–216. https://doi.org/10.1038/nrg2063

1087

1088

# Supplementary Material

# **Appendix:** Analyses assuming LD-based recombination maps

### Generating LD-based recombination rate maps

We phased variant calls using the read-aware methodology incorporated in SHAPEIT2 (Delaneau et al., 2013). For each of the mouse population samples, we carried out the following procedure. First, we created a stringently filtered set of SNPs following Booker et al. (2017), by only including biallelic variant sites that met the following criteria: no overlap with indels, no missing data, QUAL >= 30, genotype quality (GQ) greater than or equal to 15 in all individual genotypes, sequencing depth (DP) greater than or equal to 10 for all individuals, rejected sites with significant deviation from Hardy-Weinberg equilibrium at the level $p < 0.05$). Using the filtered variants, we extracted phase informative reads. We then ran SHAPEIT2 in 'assemble' mode to phase our stringently filtered variants. Finally, we converted the output of SHAPEIT2 to FASTA files, which contained two haplotypes per diploid sample using custom Python scripts.

We ran LDhelmet version 1.9 (Chan et al. 2012) on the phased haplotypes, in order to estimate the population-scaled recombination rate, $\rho = 4N_e r$, where $N_e$ is the effective population size and $r$ is the rate of crossing over between two sites per generation, for each of the mouse populations. We calculated the ancestral prior probability for each variant site that we passed to LDhelmet using the method developed by Keightley & Jackson (2018) as implemented in the program *est-sfs* v2.01. As input for this program, we generated files including each variant and invariant site that met less stringent filtering criteria than that described above (QUAL > 30, no missing data, no overlap with indels, ExcessHet < 13, no more than two alleles per site), and additionally discarded sites that did not have full outgroup information (alleles from both *Mus famulus* and *Mus pahari* for mouse samples mapped to mm10). For sites that were present in the input to LDhelmet, but not in the input for the Keightley & Jackson (2018) method, because they lacked complete outgroup data, we assigned the ancestral prior following Equation 18 in Keightley & Jackson (2018). We used the resulting information about the ancestral states of SNPs to populate the 4x4

1118    mutation transition matrix used by LDhelmet (Chan et al. 2012). To estimate fine-scale

1119    recombination rates in each of our populations, we ran the *find_confs* component of LDhelmet

1120    with a window size (-w) of 50 SNPs to generate haplotype configuration files from the phased

1121    FASTA files we made in the step above. Subsequently, we ran the *table_gen* and *pade* components

1122    of LDhelmet with the default parameters, with the exception of θ [-t] which we set to the point

1123    estimate of $\pi$ at 4-fold degenerate synonymous sites specific to each population. To estimate $\rho$ we

1124    ran the *rjmcmc* component of LDhelmet with a width [-w] of 50 SNPs, a block penalty [-b] of 100, a

1125    partition length of 4001 SNPs, an overlap of 200 SNPs, a burn-in period of 100,000 iterations

1126    followed by 1,000,000 iterations of the Markov chain.

1127

1128    Comparison of LD-based recombination rates among taxa

1129    When analysing patterns of genetic diversity under a model of selection at linked sites, the way in

1130    which recombination rate estimates were obtained may affect parameter estimates. We analysed

1131    the relationship between nucleotide diversity and genetic distance from functional elements in *M.*

1132    *m. castaneus* assuming either a high-resolution recombination map constructed using patterns of

1133    linkage disequilibrium (LD) or the pedigree-based map constructed by Cox et al. (2009). These two

1134    approaches for generating recombination rate maps have both advantages and disadvantages. By

1135    examining patterns of LD, the population-scaled recombination rate ($\rho = 4N_e r$), where $r$ is the

1136    recombination rate, can be inferred from a relatively small sample of unrelated individuals at very

1137    fine-scales. However, natural selection can influence LD and may therefore affect such

1138    recombination rate estimates (Clark et al. 2010). Alternatively, direct estimates of the

1139    recombination rate can be obtained from crossing experiments, but to achieve a high-resolution

1140    recombination map, a very large number of individuals need to be genotyped and this has typically

1141    precluded the use of whole-genome re-sequencing in some species such as mice, thereby limiting

1142    resolution.

1143

1144    We generated recombination rate maps from patterns of LD for each of the mouse taxa, and

1145    compared these to the pedigree-based estimates obtained by Cox et al. (2009). It is worth pointing

1146    out that the Cox et al. (2009) map is an estimate of the recombination map that was generated

1147    using inbred strains of mice of predominantly *M. m. domesticus* origin and there are known

1148    differences in total genetic map length and local recombination rate between *M. musculus* sub-

1149    species (Dumont & Payseur 2011). For simplicity, we treat the Cox map as a baseline comparison

1150    for each of the recombination rate landscapes we inferred.
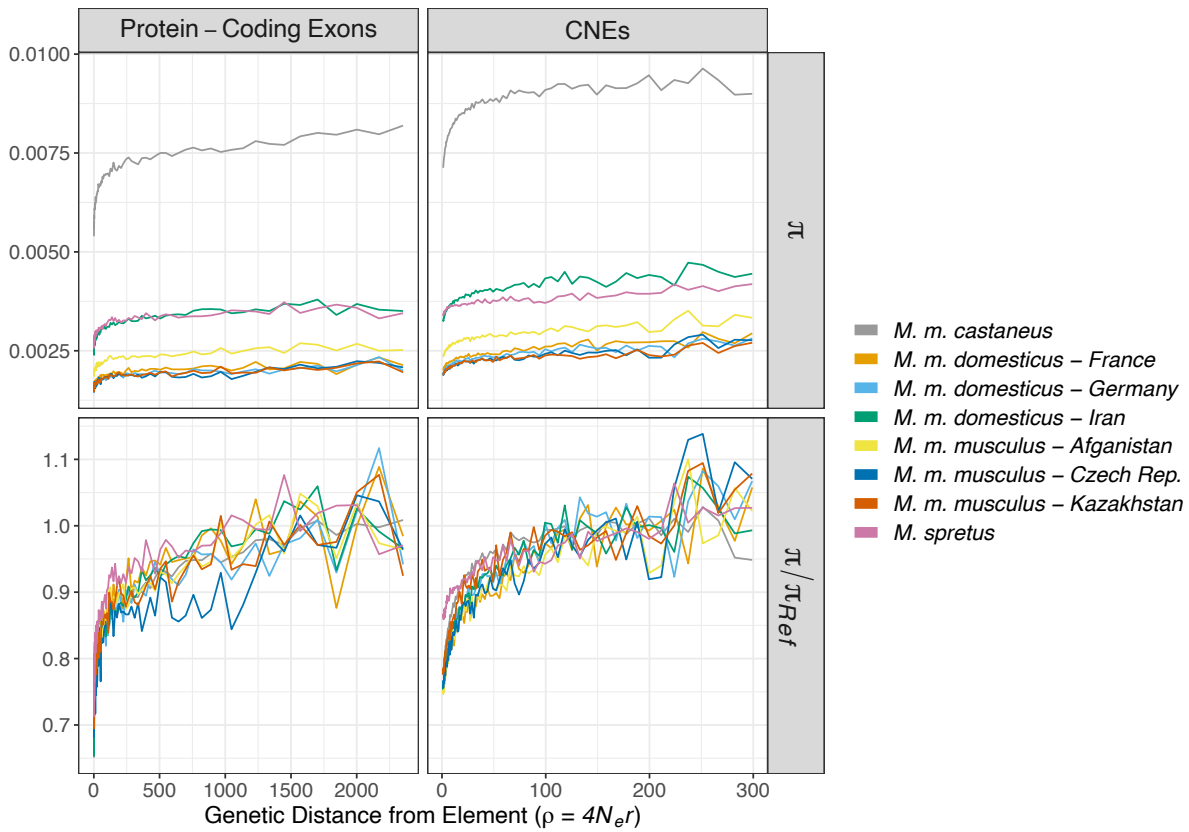
42

1151

1152    We calculated Spearman's correlation between LD-based recombination rate estimates obtained

1153    for each mouse taxa and recombination rate estimates from the Cox map in windows from 1Mbp

1154    up to 20Mbp. Across all scales tested, the recombination maps for *M. m. castaneus* and *M. m.*

1155    *musculus* from Afghanistan showed the highest level of congruence with the Cox map (Figure A.1).

1156    The correlation exhibited by the *M. m. castaneus* was very similar to the correlation previously

1157    reported (Booker et al., 2017). For the purposes of calculating genetic distances, we used the LD-

1158    based recombination rate estimates for *M. m. castaneus*.



1159

1160    **Figure A.1** Spearman rank correlation coefficients between recombination maps inferred using

1161    LDhelmet for wild mice and the pedigree-based map of Cox et al. (2009). Correlations were

1162    calculated in non-overlapping windows of discrete physical size. For the purposes of visualisation,

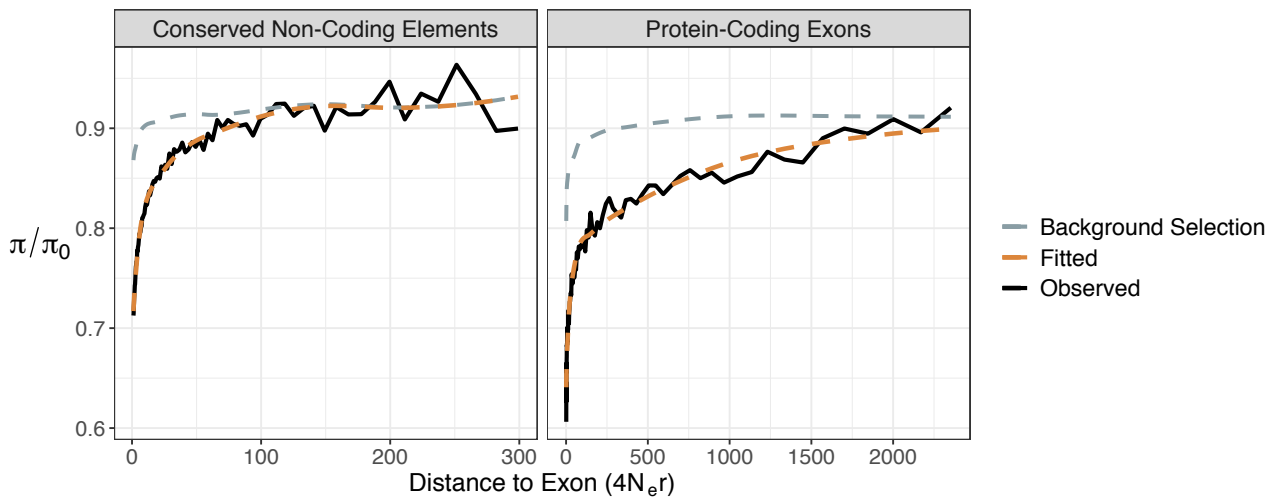1163    the confidence interval is only shown for the *M. m castaneus* map.

1164

**Figure A.2** Identical to Figure 2 in the main text except that genetic distances were calculated assuming the LD-based recombination map constructed for *M. m. castaneus*.



**Supplementary Table S1** Comparison of uSFS model fits for each taxa and class of sites considered. The maximum likelihood estimate of model parameters are shown along with the estimated uSFS. A parameter key is given as a second sheet in the spreadsheet.
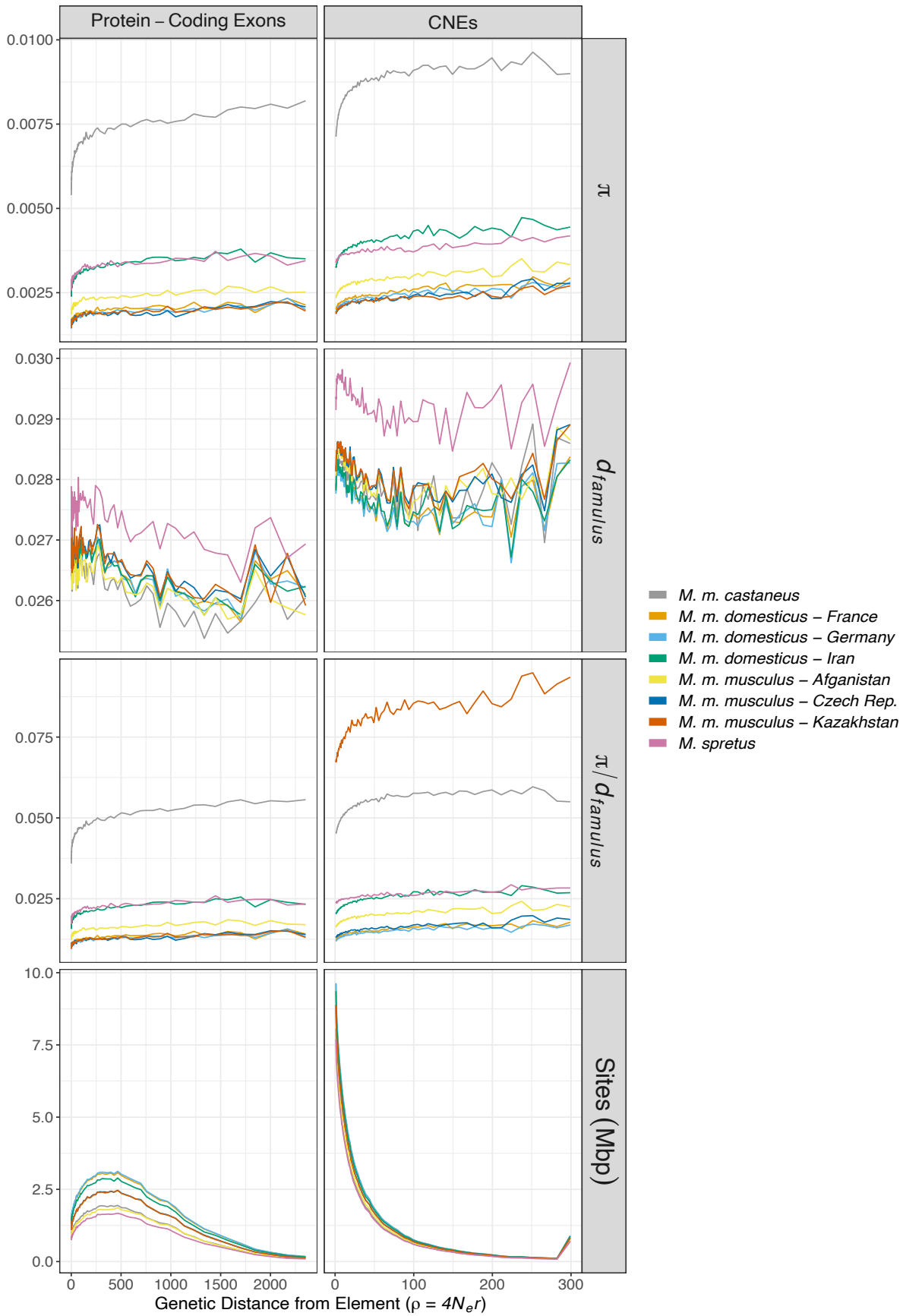
1176     **Supplementary Table S2** Parameters of the distribution of fitness effects for deleterious mutations

1177     as well as the positive selection parameters estimated for each population using *polyDFE*. Point

1178     estimates are provided as well as 95% bootstrap confidence intervals. A parameter key is given as

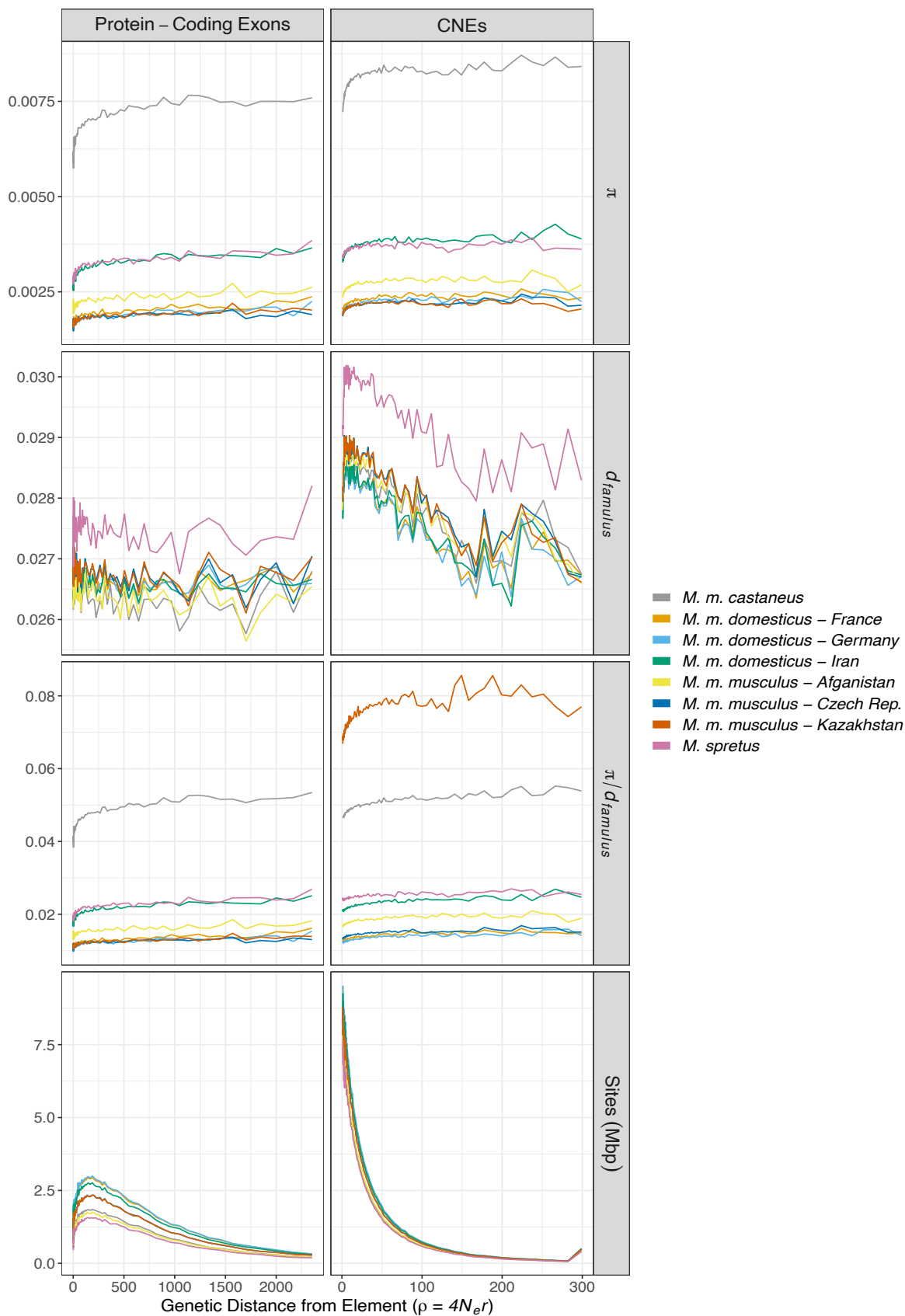1179     an additional sheet in the spreadsheet.

1180

1181     **Supplementary Table S3** Estimates of positive selection parameters obtained by fitting a models

1182     of selective sweeps and background selection to troughs in nucleotide diversity. Parameters are

1183     given for models assuming a one or two discrete classes of advantageous mutations as well as an

1184     exponential distribution of fitness effects. Estimates of the fitness change brought about by

1185     positive selection in protein-coding exons and CNEs are also given in the table. A parameter key is

1186     given as an additional sheet in the spreadsheet.

1187

1188

**Figure S1** Additional summary statistics in the regions surrounding functional elements assuming

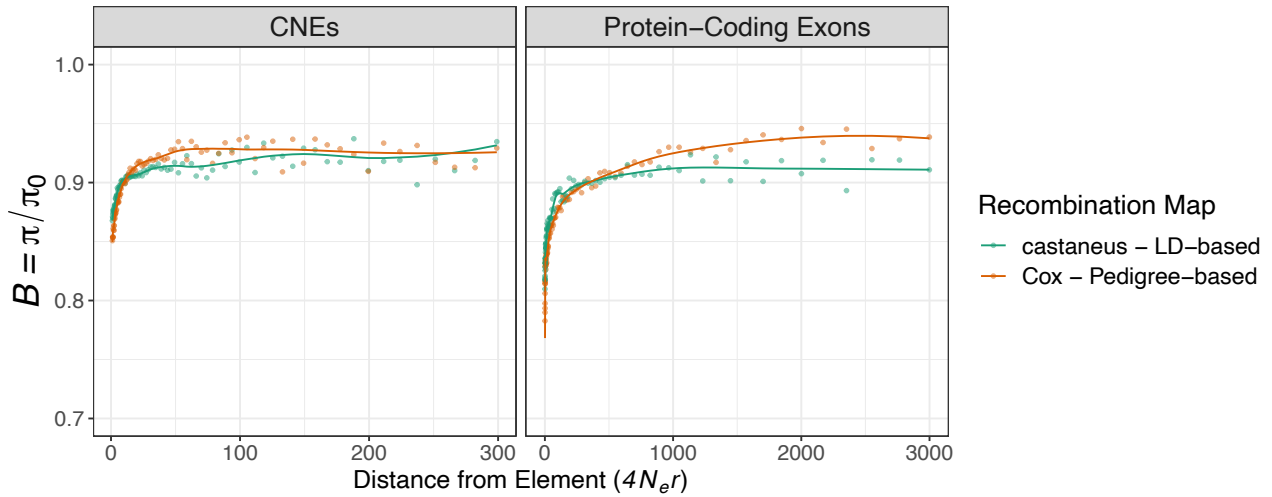the LD-based map of recombination rate variation we inferred for *M. m. castaneus.*

**Figure S2** Additional summary statistics in the regions flanking functional elements assuming the Cox map.
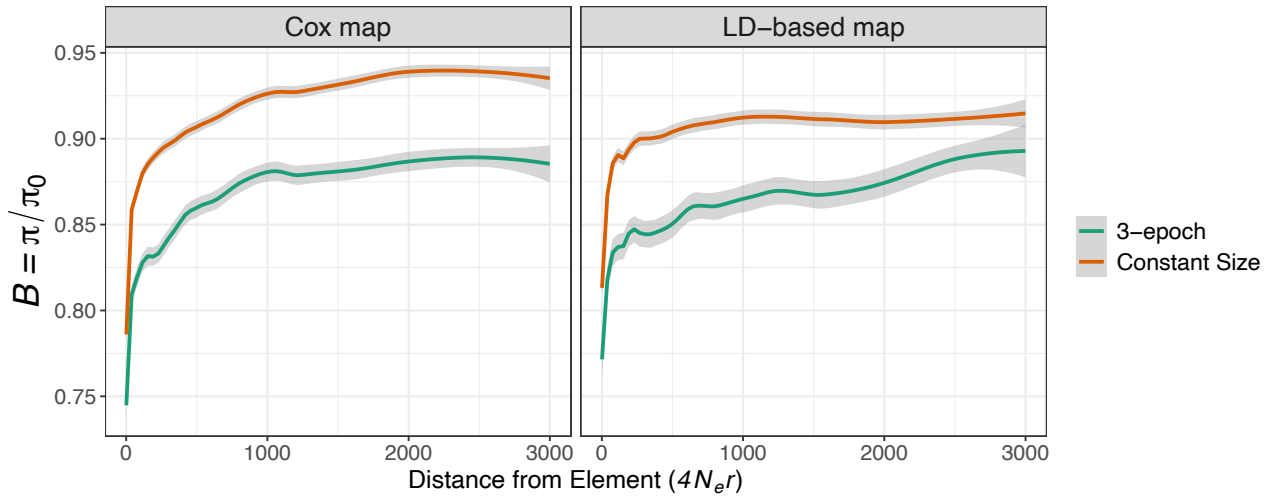
1195



1196

1197 **Figure S3** The reduction in neutral genetic diversity relative to neutral expectation caused by

1198 background selection ($B$) observed in simulated datasets. Simulations assumed either the LD-

1199 based recombination map or the pedigree-based map of Cox et al. (2009). Lines indicate the fit of

1200 a Loess regression fitted to the data with a span of 0.3 and the number of sites in each bin used as

1201 weights.

1202

1203

1204

1205



1206

1207 **Figure S4** The reductions in neutral genetic diversity relative to neutral expectation caused by

1208 background selection ($B$) observed in simulated datasets when modelling a population with

1209 constant size, or the three-epoch demographic model estimated by Booker and Keightley (2018).

1210 $\pi_0$ in the constant size simulations was 0.01. $\pi_0$ was 0.0042 in the 3-epoch simulations, which was

1211 calculated from the harmonic mean of population sizes. Lines indicate the fit of a Loess regression

1212 fitted to the data.

1213

1214