

# Large changes in detected selection signatures after a selection limit in mice bred for voluntary wheel-running behavior

## Short title: Evolution following reaching a selection limit

David A. Hillis<sup>1\*</sup>, Liran Yadgary<sup>2, #a</sup>, George M. Weinstock<sup>3</sup>, Fernando Pardo-Manuel de Villena<sup>2</sup>, Daniel Pomp<sup>2</sup>, and Theodore Garland, Jr.<sup>4</sup>

<sup>1</sup> Genetics, Genomics, and Bioinformatics Graduate Program, University of California, Riverside, California 92521

<sup>2</sup>Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina 27599

<sup>3</sup> The Jackson Laboratory for Genomic Medicine, Farmington, CT 06032 and Department of Genetics and Genome Science, University of Connecticut Health Center, Farmington, CT 06030

<sup>4</sup> Department of Evolution, Ecology, and Organismal Biology, University of California, Riverside, California 92521

<sup>#a</sup> Current Address: Hazera Seeds Ltd. in Israel, Berurim M.P Shikmim, Israel 7983700

\* Corresponding author

30 E-mail: [davidhillis@ucsb.edu](mailto:davidhillis@ucsb.edu)

31 [Abstract](#)

32 In various organisms, sequencing of selectively bred lines at apparent selection limits  
33 has demonstrated that genetic variation can remain at many loci, implying that evolution  
34 at the genetic level may continue even if the population mean phenotype remains  
35 constant. We compared selection signatures at generations 22 and 61 of the “High  
36 Runner” mouse experiment, which includes 4 replicate lines bred for voluntary wheel-  
37 running behavior (HR) and 4 non-selected control (C) lines. Previously, we reported  
38 multiple regions of differentiation between the HR and C lines, based on whole-genome  
39 sequence data for 10 mice from each line at generation 61, which was >31 generations  
40 after selection limits had been reached in all HR lines. Here, we analyzed pooled  
41 sequencing data from ~20 mice for each of the 8 lines at generation 22, around when  
42 HR lines were reaching limits. Differentiation analyses of allele frequencies at ~4.4  
43 million SNP loci used the regularized T-test and detected 258 differentiated regions with  
44 FDR = 0.01. Comparable analyses involving pooling generation 61 individual mouse  
45 genotypes into allele frequencies by line produced only 11 such regions, with almost no  
46 overlap among the largest and most statistically significant peaks between the two  
47 generations. These results implicate a sort of “genetic churn” that continues at loci  
48 relevant for running. Simulations indicate that loss of statistical power due to random  
49 genetic drift and sampling error are insufficient to explain the differences in selection  
50 signatures. The 13 differentiated regions at generation 22 with strict culling measures  
51 include 79 genes related to a wide variety of functions. Gene ontology identified  
52 pathways related to olfaction and vomeronasal pathways as being overrepresented,  
53 consistent with generation 61 analyses, despite those specific regions differing between

54 generations. Genes *Dspp* and *Rbm24* are also identified as potentially explaining  
55 known bone and skeletal muscle differences, respectively, between the linotypes.

56  
57  
58  
59

60 [Introduction](#)

61 Although evolution can result in organisms with spectacular capabilities or able to  
62 survive in exceptionally inhospitable environments, all adaptations are bound within  
63 certain limits. These limits are commonly observed in laboratory and agricultural  
64 selection experiments (Dobzhansky and Spassky 1969; Al-Murran and Roberts 1974;  
65 Careau et al. 2013; Schlötterer et al. 2015; Lillie et al. 2019). Among various possible  
66 causes of selection limits (Al-Murran and Roberts 1974; Falconer 1989; Douhard et al.  
67 2021), the simplest explanation is the loss of genetic variation, such that narrow-sense  
68 heritability declines to zero (e.g., see Brown and Bell 1961). However, selection  
69 experiments have frequently found that genetic variation remains after reaching a  
70 selection limit (e.g., Lerner and Dempster 1951; Roberts 1966; Dobzhansky and  
71 Spassky 1969; Bult and Lynch 2000; Burke et al. 2010; Careau et al. 2013; Lillie et al.  
72 2019; Hillis et al. 2020). Even for alleles favored by selection, fixation is far from  
73 guaranteed (Burke et al. 2010; Schlötterer et al. 2015; Stephan 2016; Hillis et al. 2020).

74 One selection experiment that has continued selection long after reaching a limit  
75 is the High Runner (HR) mouse experiment, which started in 1993 with the purchase of  
76 224 outbred ICR mice from Harlan Sprague Dawley (Swallow, Carter, et al. 1998).  
77 These were randomly bred for two generations, then split into ten breeding pairs to  
78 found each of eight closed lines. Four of these lines were designated to serve as non-  
79 selected control lines, while the other four were selected based on voluntary wheel  
80 running. In selected lines, all mice are given access to wheels for 6 days and the male  
81 and female of each family with the highest running on days 5 and 6 would be used as  
82 breeders (no sib-mating). After about 22 generations of selection, three of the four HR

83 lines (with the fourth line following suit a few generations later) had plateaued in their  
84 running at approximately 2.5 to 3 times as many revolutions as the controls (Careau et  
85 al. 2013). Recently, the experiment has reached its 100<sup>th</sup> generation since selection  
86 began and, with exception of some generations when the experiment moved from  
87 Wisconsin to California (generations 32 to 35) and during Covid-19 lockdowns  
88 (generations 91 to 98), selection has continued nearly uninterrupted in the interim.  
89 Whether selection interruption following the move to California resulted in statistically  
90 significant changes to running behavior has not yet been analyzed.

91 Numerous physiological and morphological differences between the HR and  
92 control lines have been documented (Rhodes et al. 2005; Swallow et al. 2009; Garland,  
93 Jr., Schutz, et al. 2011; Wallace and Garland, Jr. 2016; Khan et al. 2024). These  
94 include traits associated with motivation to run, such as changes in dopamine (Rhodes  
95 et al. 2001; Mathes et al. 2010), serotonin (Waters et al. 2013), and endocannabinoid  
96 signaling (Thompson et al. 2017), as well as changes in brain size and structure (Kolb,  
97 Rezende, et al. 2013). Additionally, changes associated with ability to run have been  
98 found, including endurance capacity (Meek et al. 2009), maximal aerobic capacity  
99 (VO<sub>2Max</sub>) (e.g., Swallow, Garland, Jr., et al. 1998; Kolb et al. 2010; Dlugosz et al. 2013;  
100 Hiramatsu et al. 2017; Cadney et al. 2021), heart size (Kolb et al. 2010; Kolb, Kelly, et  
101 al. 2013; Kelly et al. 2017), skeletal muscle physiology (Dumke et al. 2001; Syme et al.  
102 2005; Guderley et al. 2008; Castro et al. 2022), and bone morphology (Garland, Jr. and  
103 Freeman 2005; Kelly et al. 2006; Middleton et al. 2008; Middleton et al. 2010; Wallace  
104 et al. 2010; Wallace et al. 2012; Castro and Garland, Jr. 2018; Copes et al. 2018;  
105 Schwartz et al. 2018).

106        Previously, whole-genome differentiation analyses using individual mouse data  
107    from 10 males from each of the eight lines at generation 61 identified at least 13  
108    genomic regions differentiated between the control and HR lines (Hillis et al. 2020; Hillis  
109    and Garland, Jr. 2022). Within these regions were genes associated with development  
110    of the brain, heart, bones, and limbs, in addition to reward pathways, and even the  
111    vomeronasal system (see also Nguyen et al. 2020). Dropping individual lines from  
112    analyses revealed new potential signatures of selection and demonstrated that the HR  
113    lines have evolved in different ways at the genomic level (“multiple solutions” Garland,  
114    Jr., Kelly, et al. 2011) that increase wheel-running behavior (Hillis and Garland, Jr.  
115    2022). Despite being ~30-35 generations past the selection limit, a great deal of genetic  
116    diversity remained in all 8 lines including many regions identified as differentiated  
117    between the HR lines and controls.

118        With the selection limit achieved near generation 22, one might expect many if  
119    not most biologically relevant SNPs to already be differentiated by that generation.  
120    Thus, with respect to the ability to detect selection signatures, little advantage would be  
121    gained from allowing ~30-35 generations to pass before testing for allelic differentiation  
122    between the HR and control lines. Furthermore, simulations performed by Baldwin-  
123    Brown et al. (2014) demonstrate that increasing the number of generations could  
124    reduce power to detect some loci under selection, which they attributed to noise created  
125    by random genetic drift. Reasonably, one might expect that drift over enough  
126    generations may cause control lines to diverge from each other in allele frequencies,  
127    such that selection signatures are obscured in statistical tests that compare replicate  
128    sets of selected and control lines. For example, if some control lines become fixed for

129 one allele and the remaining control lines become fixed for another, then, even if all HR  
130 lines were fixed for the same allele favored by selection, statistically significant  
131 differentiation would be difficult to detect. Therefore, analyses of a generation close to  
132 when a selection limit is first reached would be optimal for tests of genetic  
133 differentiation.

134 In the present study, we analyze pooled sequence data from each of the four HR  
135 lines and four control lines at generation 22. Although these analyses identify many  
136 regions containing genes associated with systems known to be phenotypically  
137 differentiated between the HR and control linetypes, they largely differ from those  
138 previously identified with the generation 61 individual mouse sequence data (Hillis et al.  
139 2020). Furthermore, the number of differentiated regions detected at generation 22 are  
140 more than 20-fold greater than those detected with generation 61 data (treated as  
141 pooled data).

142 We first discuss possible methodological causes of these differences (e.g.,  
143 pooled vs individual mouse data) and find them lacking. We therefore develop a simple  
144 simulation model, with leptokurtic distribution of locus effect sizes, to test the possibility  
145 that a hypothetical physiological constraint on wheel running could contribute to the  
146 differences between generations 22 and 61 selection signatures. Ignoring locus effect  
147 size, results demonstrate that such constraints can contribute to a reduction in power  
148 and increased variability in the detected response to selection in generations after the  
149 selection limit. However, the magnitude of these effects appears insufficient to explain  
150 the differences observed between generations 22 and 61 in the real data. In addition,  
151 effect size was an important determinant of the ability to detect selection signatures in

152 the simulations, including a more than 2-fold increase in power to detect loci with large  
153 effect size at generation 22 as compared to generation 61. Thus, with strict culling  
154 procedures, we suspect that many of the selection signatures detected at both  
155 generations are likely to represent loci with relatively large effects on wheel running.  
156 The regions detected at generation 22 include genes related to olfactory/vomeronasal  
157 systems, which are also identified at generation 61 (Hillis et al. 2020; Nguyen et al.  
158 2020; Hillis and Garland, Jr. 2022).

159

160

161 Materials and methods

162

163 High runner mouse model

164 As described previously (Careau et al. 2013; Swallow, Carter, et al. 1998), 112 males  
165 and 112 females of the outbred Hsd:ICR strain were purchased from Harlan Sprague  
166 Dawley in 1993 and designated as generation -2. Mice would be randomly bred for 2  
167 generations (-2 and -1) with 2-3 generation -1 mice from each family randomly chosen  
168 to contribute to 1 of 8 different closed lines. Four of these lines were randomly picked to  
169 be “High Runner” (HR) lines, in which mice would be selected for breeding based on  
170 voluntary wheel running. The remaining 4 lines were used as Control (C) lines, without  
171 any selection. Generation 0 was the first generation where HR lines were paired based  
172 on running levels (10 males and 10 females for each line) with generation 1 the first  
173 product of selection.

174 Wheel running measurements were collected by giving mice at approximately 6-8  
175 weeks of age, access to wheels for six days. The amount of running (total revolutions)  
176 on days 5 and 6 was used as the selection criterion. Both days 5 and 6 are used for  
177 repeatability in running behavior and robustness against bad data for a single day  
178 (Careau et al. 2013). For the HR lines, the highest-running male and female from within  
179 each of 10 families were chosen as breeders (within-family selection). For the non-  
180 selected C lines, one male and one female from each of 10 families were chosen as  
181 breeders, independent of wheel running measurements. Sib-mating was disallowed in  
182 all lines (Swallow, Carter, et al. 1998).

183

184     [Genome sequencing and allele frequency determination](#)

185     Roughly 10 male and 10 female mice were taken from each line at generation 22 (Khan  
186     et al. 2024). Mice were decapitated without anesthesia because blood was being taken  
187     for a study of hormone levels (corticosterone) that can respond rapidly to additional  
188     handling or anesthesia. Subsequently, their DNA was extracted from tail tips and then  
189     pooled for determination of allele frequency for each line. This pooled DNA was  
190     sequenced with paired end pooled sequencing with Illumina HiSeq 2500 sequences  
191     were trimmed and aligned to the GRCm38/mm10 mouse genome assembly.

192     Generation 22 used trimmomatic v0.39 for trimming, BWA v0.7.17 for alignment,  
193     Samtools v1.14 for sorting and indexing, picard v2.26.11 for marking duplicates, and  
194     GATK v4.1.8.1 for calling SNPs. SNPs were filtered to keep those with read quality  
195     ("RQ")  $\geq$  20, DP  $\geq$  10, were missing either quality score, or missing the allele frequency  
196     all together, or had MAF  $>$  0.0126. Allele frequencies ("AF") were determined for  
197     generation 22 by taking the read depth of the alternate nucleotide allele (i.e., allele  
198     differing from the GRCm38/mm10 alignment) and dividing by the read depth for the  
199     locus. After all quality control methods were implemented, 4,446,523 loci remained for  
200     generation 22.

201           The generation 61 data were taken from Hillis et al. (2020). 80 male mice (10  
202     from each line) were subject to whole genome sequencing and reads were trimmed and  
203     aligned to the GRCm38/mm10 mouse genome assembly as described in Didion et al.  
204     (2016). This generated an average read depth of 12X per mouse. SNPs were filtered  
205     to keep those with genotype quality ("GQ")  $>$  5, read depth ("DP")  $>$  3, minimum allele  
206     frequency ("MAF")  $>$  0.0126 for all samples, and Mapping Quality ("MQ")  $>$  30. One of

207 the 80 mice was excluded due to likely contamination (as in Xu and Garland 2017),  
208 leaving 79 for the following analyses. SNPs not found to be present in at least two of  
209 the 79 mice were also removed from analysis. After all quality control methods were  
210 implemented, 5,932,148 loci remained for analyses. To allow comparison with the  
211 pooled sequencing data from generation 22, we calculated allele frequencies as the  
212 number of alternative alleles divided by 2 times the number of mice (i.e., 20 or 18 for  
213 HR3).

214

## 215 [Statistical analyses](#)

216 For generations 22 and 61 we used an arcsine-squared transformation (Ahrens et al.  
217 1990) of the AF. Analyses were conducted on both generations using a traditional T-  
218 test, regularized T-test (RegT)(Baldwin-Brown et al. 2014, see also Baldi and Long  
219 2001), and a variant of the regularized T-test which uses a sliding window to calculate  $\bar{v}$   
220 (WRT test) (S1 File). The regularized T-test was based on a Bayesian method meant to  
221 minimize the type-I errors caused by sampling error with small sample sizes (Baldi and  
222 Long 2001; Baldwin-Brown et al. 2014), such as the 8 total lines in the HR mouse  
223 selection experiment. We performed these tests and determined the permutation-based  
224 false discovery rate (FDR) for each method (see below). For comparison, we also  
225 performed the RegT and WRT tests on loci found in both generation 22 and 61 (from  
226 pooling individual mouse genotypes) data sets along with the FDR. Since standard T-  
227 tests do not require whole genome or region variances of other loci, the p-values of loci  
228 shared between the two generations could simply be extracted from the complete  
229 original analyses.

230

231 [Permutation-based false discovery rate](#)

232 To determine relative power of generation 22 allele frequencies with arcsine-square  
233 transformation using T-test, regularized T-test, and WRT test, we attempted to calculate  
234 a critical threshold by estimating the FDR of 10% (Benjamini and Hochberg 1995; Xie et  
235 al. 2005). However, after calculating p-values for complete permutations of the different  
236 lines within linetype to better understand the null distribution, we concluded that this  
237 estimated FDR was underestimating the true false discovery rate. Therefore, using  
238 these same permutations, we calculated the FDR directly.

239 Direct calculations of FDR were performed by calculating FDR for each locus of  
240 the unpermuted data whose p-value was below 0.01 in accordance with the equation:

$$241 \quad FDR = \frac{n \text{ False Positives}}{n \text{ rejected Null Hypotheses}}$$

242

243 This was implemented for each locus with:

$$244 \quad FDR = \frac{\frac{n \text{ permuted loci significant at } p}{35}}{n \text{ unpermuted loci significant at } p}$$

245

246 Loci with nominal p-value < 0.05 were ordered by FDR score, the p-value was identified  
247 for the locus with the largest FDR below 0.01, and any p-values less than or equal to  
248 the p-value for this locus was treated as significant. The SNPs with FDR = 0.01 were  
249 then further grouped into “significant regions” by grouping any loci within 1mbp of  
250 another and separating groups whose closest SNPs are further than 1mbp.

251

252 Divergence over time  
253 To test for a difference in the number of loci showing a significant change in allele  
254 frequency between the HR and C lines from generations 22 to 61, we first conducted a  
255 paired T-test for the 4 C lines and separately for the 4 HR lines. These tests were  
256 based on eight data points for each linetype, i.e., the mean allele frequency for each line  
257 at a given locus at generation 22 and 61. The T-score for the C T-test was then  
258 subtracted from the T-score for the HR T-test, and the absolute value was taken. This  
259 was repeated for each locus, producing values for approximately 2 million loci  
260 (excluding where either the C or HR T-test failed for numerical reasons). These  
261 analyses were then repeated with all 35 permutations (as described above) to estimate  
262 the null distribution of the score based on  $\sim 2,000,000 * 35 = \sim 75,000,000$  values.  
263 These scores were ordered to identify the 5<sup>th</sup> percentile threshold for comparison with  
264 the distribution of the unpermuted results.

265  
266 “Strict” culling for biological and AF change analyses  
267 Rather than attempt to focus on the genes of more than 100 regions for each of the  
268 different statistical tests, analyses of biological significance and comparisons of change  
269 in allele frequencies between generations 22 and 61 were done using a subset of the  
270 regions identified by FDR. WRT and regularized T-test first culled by removing regions  
271 containing only one significant locus, then culled such that only regions containing at  
272 least 20 significant loci or the lowest p-value among loci in the region was below 1.00E-  
273 04. Regions associated with the T-tests were culled in a similar manner as the WRT  
274 and RegT test, except the p-value cutoff used was 1.00E-06 due to naturally lower p-

275 values. These culling methods should also serve to reduce the influence of sampling  
276 error, as it would be increasingly unlikely for sampling error to simultaneously  
277 underestimate among-line variance across multiple linked SNPs and lines. We will refer  
278 to these additional culling methods below as “strict” culling.

279

280 [Comparison of selection signatures in generations 22 and 61](#)

281 Changes in allele frequencies from generation 22 to 61 were analyzed for each region  
282 identified by strict culling for generations 22 and 61. For regions significant at  
283 generation 22, each region and its included SNPs with nominal  $p < 0.05$  at generation 22  
284 were matched with SNPs at generation 61. The allele frequencies of these SNPs were  
285 averaged for each line and generation and line graphs created (one for each line) with  
286 generation 22 AF on the left and generation 61 AF on the right. This was then repeated  
287 for regions significant at generation 61, except each region and its included SNPs with  
288 nominal  $p < 0.05$  at generation 61 were matched with SNPs at generation 22.

289

290 [Simulations to compare presumptive statistical power across generations](#)

291 The available data from the two generations differ in multiple ways that might affect  
292 cross-generation comparisons of selection signatures. Each generation, each line is  
293 reduced to ~20 individuals when ~10 breeding pairs are formed. An ideal “sample” from  
294 a given generation would include all 20 of those breeding individuals. Instead, our  
295 sample from generation 22 was of ~10 males and 10 females per line that were  
296 sampled at random at the time of weaning (i.e., they were not the 20 breeding parents).  
297 In contrast, the mice from generation 61 were a semi-random sample of 10 males from

298 each line (except nine from HR3 and one female that was unintentionally used from  
299 another line) (Hillis et al. 2020).

300 For a pooled DNA sample, as for generation 22, a further ideal condition is for  
301 the sample of DNA from each mouse to be of equal volume and concentration through  
302 the extraction and pipetting steps prior to pooling. This would then result in each  
303 mouse's alleles being represented in equal quantities in the pooled sequencing sample.

304 The next source of error is read depth, which is effectively a random sampling of  
305 alleles from the pooled sample. Our generation 22 samples were read at an average  
306 depth of 24X. Thus, the frequency of alternative nucleotide alleles for a given SNP  
307 locus was calculated by counting the number of alternative alleles, which was taken as  
308 anything other than the reference. Thus, not all of the 40 alleles (as one of two  
309 possibilities) contributed by the 20 mice could have been identified with a read depth of  
310 24X, which acts as 24 samples taken with replacement.

311 The generation 61 data are from individual sequencing of 10 mice per line at an  
312 average read depth of 12X, with those results then used to predict the genotype for  
313 each SNP and mouse (Hillis et al. 2020). This should allow for the representation of  
314 nearly all alleles ( $N = 2$  alleles  $\times$  10 mice). Originally, those data were analyzed as such  
315 via mixed models to detect selection signatures (Hillis et al. 2020). Here, to allow  
316 comparison with the pooled sequencing data from generation 22, we calculated allele  
317 frequencies as the number of alternative alleles divided by 2 times the number of mice  
318 (i.e., 20 or 18 for HR3), which should incorporate 19-20 unique alleles in equal  
319 proportion. Given that the data available from the two generations differ in multiple  
320 ways, we used simulations in an attempt to assess how this might affect our results.

321 For generation 22, simulations to elucidate possible sampling errors were  
322 performed such that alleles for 20 mice were sampled using a random binomial  
323 distribution assuming population allele frequencies of (0.05, 0.10, 0.15, ..., 0.90, and  
324 0.95). Then an allele depth was randomly sampled from the actual quality data for the  
325 SNPs used in the generation 22 analyses and alleles were sampled from these  
326 simulated 20 mice (with replace) equal to this read depth. The allele frequency was  
327 then calculated as the number of alternative alleles (1) divided by the total read depth.  
328 This generated a distribution of allele frequencies given a particular starting AF for the  
329 population and was repeated 100,000 times for each starting population AF.

330 For generation 61, simulations were performed such that alleles for 10 mice  
331 were sampled using a random binomial distribution assuming population allele  
332 frequencies of (0.05, 0.10, 0.15, ..., 0.90, and 0.95). Then for each simulated mouse's  
333 genotype, a genotype quality was randomly sampled from the actual quality data for the  
334 SNPs used in the generation 61 analyses. If the simulated genotype for the mouse was  
335 heterozygous, then the genotype quality would be used to generate a 0 or 1 with the  
336 probability of a 1 equaling that of the probability of a genotyping error. If a 1 was  
337 generated (thus an error occurred) the second allele for the mouse was replace with a  
338 copy of the first allele of the mouse. The allele frequency was then calculated as the  
339 number of alternative alleles (i.e., 1) for all ten mice divided by the total alleles (i.e., 20).  
340 This generated a distribution of allele frequencies given a particular starting AF for the  
341 population and was repeated 100,000 times for each starting population AF.

342 Power analyses were then done by sampling four AF values from the simulated  
343 AF values from an actual population AF of 0.4 for one linetype. Likewise, four AF

344 values were sampled from the simulated AF values from an actual population AF of 0.6  
345 for the other linetype. Sampled allele frequencies were transformed using an arcsine-  
346 squared transformation. A T-test (assuming unequal variance) was then conducted  
347 comparing these 8 sampled AF values. Note that this could not be done for RegT and  
348 WRT tests because it would require simulations of regional or genome-wide variance  
349 structure. These sampling and T-tests were repeated 10,000 times.

350

351 Simulations comparing power with and without a biological constraint  
352 We used simulations to begin to address whether a biological constraint on a trait under  
353 selection (e.g., wheel running) might affect (1) the ability to detect selection signatures  
354 at generations before (e.g., generation 22) versus long after (e.g., generation 61)  
355 selection limits were reached, (2) the consistency of those signatures across  
356 generations, and (3) the rate at which loci with different allelic effect sizes respond to  
357 selection. Our rationale for using a constraint model is explained in the Discussion. As  
358 a heuristic, some of the parameters in these simulations were chosen to approximate  
359 values observed in the selection experiment and help build a model of architecture for  
360 wheel running in the HR and control mice (Sella and Barton 2019).

361 Running levels were calculated based on the general equation:

$$y = \mu + v_g + v_e$$

363 Where  $y$  is equal to the phenotype (wheel revolutions/day) of an individual mouse;  $\mu$  is  
364 the "base" mean number of revolutions (held constant at the starting value set at  
365 generation 0);  $v_g$  is the variance contributed by genetic variation; and  $v_e$  is the variance  
366 contributed by environmental effects.

367

368 As a regression model, this equation is:

369

370  $y = \mu + \beta_1 X_1 + \beta_2 X_2$

371 where the genetic variance is represented by  $\beta_1 X_1$  and the environmental variance is  
372 represented by  $\beta_2 X_2$ .  $X_1$  represents the summed effect on wheel running of all alleles  
373 carried by the individual, where, to simulate a leptokurtic distribution (Barton and Turelli  
374 1989; Reeve 2000; Reeve and Fairbairn 2001), these alleles are coded as having  
375 variable allelic affects (specifically,  $\pm 0.4, \pm 0.8, \pm 1.6, \pm 3.2, \dots \pm 204.8$ ) at frequencies  
376 inversely proportional to their effect size (specifically, 720 loci with effect  $\pm 0.4$ , 480 loci  
377 with effect  $\pm 0.8, \dots 8$  loci with effect  $\pm 204.8$ ) for a total of 2,096 loci, which approximates  
378 the number of haplotype blocks observed across all eight lines (Hillis et al. 2020).  $X_2$   
379 provides the random element of the environmental variance and is determined by  
380 randomly sampling from a normal distribution with mean = 0 and SD = 1.

381

382 The equation we applied for these simulations is:

383  $y = 4,570 + 1.3X_1 + 2,100X_2$

384

385 The values for  $\beta_1$  (1.3),  $\beta_2$  (2,100) and for the number of loci were determined in  
386 conjunction with one another to approximate realistic (in no particular order) (1)  
387 heritability of wheel running at the base generation being about 0.32 (Careau et al.  
388 2013), (2) within-line coefficients of variation as being about 0.57 (Swallow, Carter, et al.

389 1998), and (3) realistic response to selection in the HR lines (i.e., achieving ~16,000  
390 revolutions around generation 22)(Careau et al. 2013).

391  
392 However, this equation does not adequately simulate seasonal variation (see Appendix  
393 S5 in Careau et al. 2013), so we applied an additional modifier:

394 
$$y = S * (4,570 + 1.3X_1 + 2,100X_2)$$

395  
396 S is a constant that alternates cycles between 0.769 (summer), 1 (winter and fall), and  
397 1.3 (winter). As generation time for the first 61 generations was consistently around 3  
398 months, these constants can alternate with each generation. The mean of 4,570  
399 (revolutions/day) was picked to approximate the empirically determined starting running  
400 levels at generation 0 (Swallow, Carter, et al. 1998).

401 Any running level calculated as below 100 was set to 100, which is approximately  
402 the lowest amount of running ever observed. The maximum wheel-running for  
403 unconstrained simulations was 50,000 revolutions, which is nearly twice as high as has  
404 ever been observed in actual measurements from the selection experiment (Rhodes et  
405 al. 2003; Careau et al. 2013). In practice, the highest running level produced by the  
406 unconstrained simulations was 38,875 (of 24,400 total individuals simulated over 61  
407 generations for the HR lines).

408 For the starting population of any given line, two alleles were first assigned to  
409 each of the 2,096 independently segregating starting loci for 20 mice (based on the  
410 actual selection procedures: Swallow, Carter, et al. 1998) using a random binomial  
411 distribution with  $p = 0.5$ . For control lines, mice were paired, and alleles sampled from

412 each of the pair to produce two male and two female offspring (to match the number of  
413 mice that are typically retained and wheel-tested in the selection experiment). The first  
414 of each sex for each family was then chosen to contribute to the next generation, which  
415 is functionally equivalent to the selection experiment, where breeders are chosen a  
416 random within family and sex for control lines. For HR lines, alleles were sampled from  
417 the parents for each of five males and five females (typical litter size is 10). Running  
418 distances were then calculated for all offspring, and the male and female with the  
419 highest running levels within each family were selected to breed for the subsequent  
420 generation (again, based on the actual selection procedure, which uses within-family  
421 selection). For both linetypes, siblings were barred from pairing (following the selection  
422 experiment). Simulations were run for 61 generations and alleles for all breeding pairs  
423 were saved at generation 0 and every 5 generations through 60, as well as generations  
424 22 and 61. This was then repeated for 4 control lines and 4 HR lines.

425 We modeled the constraint on wheel running as a trait that itself can evolve. To  
426 obtain a realistic value for the constraint, we applied the same principles as for the  
427 wheel-running equation:

$$428 \quad C = S * (10,000 + 1.0X_{C1} + 1,750X_{C2})$$

429  
430 C is the constraint to be applied to the mouse's wheel running. S is the same seasonal  
431 multiplier used in the wheel-running equation, without which, higher running levels in  
432 winter become truncated.  $X_{C1}$  represents the genetic component of the constraint  
433 determined by (arbitrarily) 100 loci with effect sizes of  $\pm 1$  (N=48),  $\pm 4$  (N=24),  $\pm 11$   
434 (N=12),  $\pm 36$  (N=8),  $\pm 101$  (N=5), and  $\pm 306$  (N=3).  $X_{C1}$  represents the environmental

435 component, determined by sampling from a normal distribution with mean = 0 and SD =  
436 1 (similarly to wheel running). These values result in a narrow-sense heritability of ~0.2.  
437 Despite targeting a wheel-running constraint of about 16,000 revolutions in the HR lines,  
438 the base constraint value is set to 10,000 because the alleles that increase constraint  
439 are favored by selection in the HR lines. Thus, a lower base value is needed for HR  
440 lines to stop responding to selection at about 16,000 revolutions. For the constrained  
441 simulation, if a mouse ran more than its determined constraint then its revolutions were  
442 treated as equal to the constraint itself before picking the breeders for the next  
443 generation.

444 These simulations were repeated 100 times (with 4 HR lines and 4 control lines  
445 in each simulation) assuming no constraint and 100 times with the constraint (see S2  
446 File, for parameters). T-tests assuming unequal variance between the 4 control lines  
447 and the 4 HR lines were performed at each of these “saved” generations (0, 5, 10, etc.)  
448 for the allele frequencies at each locus, with an arcsine-squared transform (Ahrens et al.  
449 1990). Power was then calculated for each simulation at each saved generation by  
450 dividing the number of loci with  $p \leq 0.05$  by the total number of loci ( $N = 2,096$ ). Power  
451 was also calculated separating loci by effect size (see below).

452 Standardized selection differentials were calculated following Careau et al.  
453 (2013), by subtracting from the mean running for each sex and family the running level  
454 of the bred individual from that litter and dividing the difference by the standard  
455 deviation of the sex for that litter. Relative power under the constrained and  
456 unconstrained models was calculated using unpaired T-tests (unequal variance) on the  
457 previously described power calculations for each simulation and for each saved

458 generation. Relative power across generations was also calculated using unpaired T-  
459 tests (unequal variance), separately for constrained and unconstrained simulations.  
460 Relative consistency in detected selection signatures was calculated by first identifying  
461 the specific significant loci (at a nominal  $\alpha = 0.05$ ) at generations 22 and 61 in each  
462 simulation. Then, the percentage of loci found significant at generation 22 that  
463 remained significant at generation 61 was calculated. Unpaired T-tests (unequal  
464 variance) were performed comparing these percentages for the constrained simulations  
465 versus the unconstrained simulations. Lastly, ability to detect loci with different effect  
466 sizes was compared using a T-test (unequal variance) of the number of significant loci  
467 ( $p \leq 0.05$ ) identified for each effect size and each simulation for generation 22  
468 constrained vs unconstrained models, generation 61 constrained vs unconstrained  
469 models, constrained generation 22 vs generation 61, and unconstrained generation 22  
470 vs generation 61. For all graphs and estimates that required the calculation of a mean  
471 value, missing values were excluded from the calculations. For example, if a p-value  
472 could not be calculated for a given locus due to fixation across all lines for the same  
473 allele, then this locus would be excluded from the power analyses.

474 Analyses were performed again implementing possible sampling error calculated  
475 by the simulations to compare statistical power, as described in the previous section  
476 "Simulations to Compare Presumptive Statistical Power Across Generations". This was  
477 implemented by taking the actual allele frequency for each line at generations 22 and 61  
478 in the simulations using the constraint model. These allele frequencies were then  
479 replaced with an allele frequency sampled from the results (rounded to the nearest  
480 0.05) of the population allele frequency of the sampling error simulations (i.e., 0.05,

481 0.10, 0.15... 0.95). For example, if the allele frequency for a given line at generation 22  
482 (constraint model) was 0.25, then this 0.25 would be replaced by a randomly sampled  
483 estimated allele frequency from the sampling error simulations (generation 22) where  
484 0.25 was the actual population allele frequency. Generation 61 allele frequencies were  
485 similarly replaced using the results of the generation 61 sampling error simulations.

486

487 [Ethics statement](#)

488 The selection experiment has been carried out in strict accordance with the approval  
489 from the Institutional Animal Care and Use Committee (IACUC) at two different  
490 institutions and under multiple protocol number. All experiments have been conducted  
491 to minimize distress to the animals. Any injuries or illness were treated in accordance  
492 with veterinarian recommendations. The present manuscript uses only published  
493 sequence data and new sequence data from historical tissue samples.

494 Results

495 Basic characteristics of genetic variation

496 The number of variable loci used in the present study includes 4,446,523 for generation  
497 22 and 5,932,148 for generation 61. Generation 61 data had an average read depth of  
498 12X per mouse for 10 mice in each of the 8 lines, producing an average read depth of  
499 over 100 per line for detection of many more variable SNPs in each line. The overlap of  
500 base positions between generations 22 and 61 was 2,045,546 SNPs. As expected,  
501 minor allele frequency (MAF), generally decreases for both HR and C lines between  
502 generations 22 and 61 (Fig 1). MAF values for HR and C lines are generally similar at  
503 generation 22; however, these diverge for many regions by generation 61.

504

505 **Fig 1. Average minor allele frequencies.** Average minor allele frequencies for  
506 generation 22 control lines, generation 22 HR lines, generation 61 control lines, and  
507 generation 61 HR lines by chromosome (numbered on the right). Regions identified as  
508 differentiated at generation 22 are indicated with an orange line above each  
509 chromosome's graph (regions smaller than 50 kbp are omitted). Regions identified as  
510 consistently differentiated at generation 61 (Hillis et al., 2020) are indicated similarly  
511 with a green line.

512

513 Differentiated SNPs and chromosomal regions

514 For analyses containing all generation 22 loci (N = 4,446,523), WRT identified 1,184  
515 differentiated loci based on 0.01 FDR (Table 1). These loci fall into 258 unique regions  
516 (separated by at least 1 million base pairs). At generation 61, 1,449 loci were identified

517 as differentiated based on 0.01 FDR. Although identifying similar numbers of loci as the  
518 generation 22 analyses, P-values for individual SNPs for generations 22 and 61 show  
519 little similarity (Figs 2B-C), with arcsine-square transform Pearson's  $r = 0.116$ .  
520 Ultimately, the SNPs identified at generations 22 and 61 were largely different.  
521 Moreover, the SNPs identified at generation 61 clustered into only 11 unique regions, as  
522 compared with the 258 regions for generation 22 (Fig 3A-B).

523

524 **Fig 2. P-value comparisons between generations.** Scatterplot comparisons of the  
525 generations 22 and 61 p-values with Pearson's correlation: (A) Generation 22  
526 regularized T-test vs generation 22 WRT test ( $\text{cor} = 0.9997$ ). (B) Generation 22 WRT  
527 test vs generation 61 WRT test ( $\text{cor} = 0.0909$ ). (C) Generation 22 WRT test vs  
528 generation 61 WRT test ( $\text{cor} = 0.1156$ ). (D) Distribution of raw p-values (generation 22).

529

530 **Fig 3. Manhattan plots and volcano plots of differentiation analyses.** Manhattan  
531 plots for results from (A) the generation 22 WRT test (shared loci), (B) the generation 61  
532 (pooled) WRT test (shared loci). The red peaks indicates those that exceeded critical  
533 threshold (FDR = 0.01) for that individual test. Volcano plot including -logP vs HR allele  
534 frequency minus C allele frequency for (C) generation 22 and (D) Generation 61  
535 (orange points indicate HR AF > C AF; green points indicate HR AF < C AF). (E)  
536 Scatterplot comparing the -log p-values of the generation 61 mixed model analyses  
537 (individual mouse) with -log p-values produced when these same data are treated as  
538 pooled sequencing allele frequencies and analyzed with WRT test. Red line has  
539 intercept = 0 and slope = 1. Green line represents the least squares regression line.

540

541 **Table 1. WRT results.**

Data	Total Loci	FDR 0.01 (-logP)	Significant SNPs	All Regions	Regions after strict culling
Gen22AF	2,045,546	2.66	630	187	6
Gen61AF	2,045,546	3.06	1,285	11	4
Gen22AF	4,446,523	2.62	1,184	258	13
Gen61AF	5,932,148	3.23	1,449	11	5

542 Number of SNPs listed represents those that are statistically significant based on a False  
 543 Discovery Rate of 1% using permutations. Analyses with 2,045,546 loci incorporate only loci  
 544 which are shared between generations 22 and 61. Regions distinguished by being separated  
 545 from the next closest significant locus by more than 1 million bp. Additional regions remaining  
 546 after additional culling methods have either 20 significant loci or at least 2 significant loci with  
 547 one having a p-value <1.00E-04.

548

549 Given such notable differences between the SNPs and regions implicated by  
 550 generation 22 and 61 analyses (Table 2), analyses were repeated focusing only on the  
 551 loci found in both data sets (N = 2,045,546). With fewer loci being analyzed, fewer  
 552 significant SNPs were identified at FDR = 0.01, as well as fewer regions for all analyses  
 553 except for WRT with generation 61. The total peaks identified when using only the  
 554 shared SNPs includes 187 and 11 regions for generations 22 and 61, respectively.

555

556

**Table 2. Genomic regions identified as differentiated under “strict” culling**

G22 Region	G22 WRT	G61 Region	G61 WRT	Chr	minPOS	maxPOS	Size	Loci	Shared Loci
1	x			1	152,318,219	153,239,876	921,658	40	25
2	x			2	78,021,909	78,974,325	952,417	3	0
		1	x	3	51,199,110	51,602,693	403,584	124	65
3	x			5	32,384,612	32,975,871	591,260	32	4
4	x			5	102,846,390	106,315,986	3,469,597	63	37
		2	x	6	<b>40,933,658</b>	<b>41,748,676</b>	<b>815,019</b>	<b>5</b>	<b>1</b>
5	x			6	122,815,876	124,446,843	1,630,968	43	20
		3	x	9	<b>41,413,436</b>	<b>42,478,817</b>	<b>1,065,382</b>	<b>1,277</b>	<b>647</b>
6	x			9	80,349,989	82,894,555	2,544,567	27	20
7	x			10	14,067,617	18,376,599	4,308,983	25	7
8	x			10	20,890,526	21,419,406	528,881	33	4
		4	x	10	104,966,751	105,529,701	562,951	2	0
9	x			13	46,088,694	46,866,721	778,028	33	13
<b>10</b>	<b>x</b>			<b>14</b>	<b>52,115,206</b>	<b>53,776,455</b>	<b>1,661,250</b>	<b>42</b>	<b>7</b>
11	x			14	77,333,032	78,080,942	747,911	4	3
		5	x	15	<b>19,245,017</b>	<b>20,197,326</b>	<b>952,310</b>	<b>27</b>	<b>17</b>
12	x			18	57,707,454	60,118,834	2,411,381	128	81
13	x			18	78,018,740	78,504,680	485,941	4	4

557

A test is deemed to have produced a differentiated region if the region contains at least 20 SNPs significant at FDR = 0.01 or at least 2 SNPs significant at FDR = 0.01 and at least one SNP with p-value < 1.0E-04 (See Methods: Determination of Selection Signatures). “Loci” listed are those significant at FDR = 0.01 (Table 1), the counts themselves match the number of differentiated loci identified by the statistical test with the most such loci. “Shared Loci” are the number of loci listed in the “Loci” column that are also shared between both generations.

562

**Bolded** loci match “consistent” regions identified by Hillis et al. (2020)

563

564

565 [Regions after “strict” culling](#)  
566 Using all available SNPs for generation 22, after applying “strict” culling (see Methods),  
567 the remaining regions were reduced to 13. All of the regions implicated by these  
568 analyses included or were near genes with intuitive implications for running behavior  
569 (see Discussion). For generation 61, strict culling reduced the total peaks to only 5  
570 unique regions.

571       Despite the HR lines reaching selection limits around generation 22 or shortly  
572 thereafter (Careau et al. 2013), the most differentiated 13 regions (Table 2) have little  
573 fixation. Of the SNPs in these regions ( $N = 79,198$ ), only about 8.78% are fixed in the  
574 HR lines, which is not significantly different from the 9.21% fixed in the control lines  
575 (unequal variance t-test comparing % fixed in the 4 HR versus 4 C lines: p-value =  
576 0.4322). If we repeat this fixation comparison for the loci shared between generations  
577 22 and 61 ( $N = 42,745$ ), 1.62% are fixed in the HR lines, which is still not significantly  
578 different from the 1.58% fixed in the control lines (p-value = 0.6129).

579  
580 [Comparison of selection signatures at generation 61 for individual vs. pooled](#)  
581 sequencing data  
582 Originally, the generation 61 individual mouse data were analyzed using mixed models  
583 (Hillis et al. 2020). We compared the previously published p-values from those  
584 analyses with the p-values produced after pooling data by line and analyzing by the  
585 WRT test (Fig 3e). The mixed model analyses produced lower p-values in general, as  
586 would be expected due to loss of power with pooling (Xu and Garland 2017), with the  
587 difference being greater for lower p-values. As a result, fewer SNP loci and hence

588 fewer chromosomal regions were identified as significantly differentiated between the  
589 HR and C lines with pooled data. Of the total regions detected with FDR = 0.01, 7 were  
590 identified at generation 22 that matched the 13 “consistent” regions identified with the  
591 mixed model analyses (Hillis et al. 2020). The 6 consistent regions that were not  
592 identified by analyses of the pooled data tended to have relatively large p-values for  
593 individual SNP loci or cover a narrower area of the genome, as compared with the other  
594 7 consistent regions.

595

#### 596 [Divergence over time](#)

597 The 5th percentile threshold for the difference in T-scores determined by permutations  
598 was 5.139787. About 6.44% of the T-scores for unpermuted data were larger than this.  
599 This difference of 6.44 - 5% indicates that ~1.44% of our values for the real data  
600 (approximately 28,470 SNPs) may be considered nominally statistically significant for  $\alpha$   
601 = 0.05. This result provides statistical support for our claim that the selection signatures  
602 differ between generations 22 and 61 (see Discussion). Defining a region as containing  
603 at least 20 significant SNPs with no adjacent SNPs separated by more than 1mbp, and  
604 considering the 1,400 most significant SNPs, they cluster into 14 regions on 13  
605 chromosomes (Table 3).

606

607

**Table 3. Genomic regions of divergent evolution**

Chr	minPOS	maxPOS	Loci	Median_T	Highest_T	T_position	g22	g61 <sup>a</sup>
1	152,255,010	153,208,066	44	48.9	94.9	152,795,939	WRT	
1	156,267,494	156,908,946	21	48.5	115.5	156,699,891		
2	153,709,397	157,032,971	20	46.0	102.8	154,908,418		1 Test
3	45,831,668	52,497,670	49	51.0	139.9	51,543,977		2 Tests
4	89,020,582	90,615,570	26	52.5	103.3	90,007,699		1 Test
5	107,675,741	111,271,760	30	55.2	232.2	109,810,077		Consistent
9	41,416,364	42,248,169	31	47.3	99.8	41,533,058		Consistent
10	101,702,890	105,671,151	154	47.4	177.0	103,108,543		2 Tests
12	109,050,203	110,779,901	28	46.6	109.2	109,228,613		
14	96,560,689	98,613,561	76	49.5	212.2	97,831,005		Consistent
15	18,635,736	20,608,793	39	51.4	208.5	19,984,048		Consistent
16	45,132,582	47,948,007	23	46.5	163.9	45,158,346		
18	69,603,969	74,277,731	62	46.7	147.5	73,016,958		
19	35,121,427	35,736,790	40	48.7	250.1	35,699,388		

608 <sup>a</sup>Tests here are the three tests used by Hillis et al. 2020 (local maxima, haplotype, and FixedHR/PolyC). “Consistent” is the term

609 used to describe regions identified by all three tests.

610

611

612 [Simulations to compare presumptive statistical power across generations](#)  
613 Simulations were conducted to gauge how much the allele frequencies determined  
614 through sequencing reflect allele frequencies of the actual populations at generations 22  
615 and 61. Generation 22 allele frequencies have greater variance from the actual  
616 population AF than generation 61 (see Figs 4A-B for an example of the 0.5 population  
617 AF distribution). The greater error variance in generation 22 is associated with reduced  
618 statistical power of 0.3864 versus 0.5031 for generation 61 when comparing simulated  
619 allele frequencies of 0.4 and 0.6 (Figs 4C-D).

620

621 **Fig 4. Variance and power simulation results.** Simulations for a population allele  
622 frequency of 0.5 for (A) generation 22 and (B) generation 61. See Methods for details.  
623 Values shown are allele frequencies for each of 100,000 simulated data sets for a single  
624 line. Methodological differences in the sampling of mice and sequencing procedures for  
625 the two generations result in greater sampling error for generation 22 (i.e., larger SD).  
626 Note that binning is done such that loci that fall on a break (e.g., 0.05) are grouped into  
627 the lower bin (e.g., 0 to 0.05). Similar simulations were then conducted to create data  
628 sets for use in estimating statistical power for detecting selection signatures for  
629 generations 22 and 61. (C) Distribution of p-values for simulated allele frequencies of  
630 0.4 versus 0.6, for generation 22. Power for  $\alpha = 0.05$  is 0.3864. (D) Distribution of p-  
631 values for simulated allele frequencies of 0.4 versus 0.6, for generation 61. Power for  $\alpha$   
632 = 0.05 is 0.5031.

633

634 [Simulations comparing power with and without a biological constraint](#)  
635 Simulations were performed modeling response to selection assuming either a  
636 constraint with a base of 10,000 revolutions per day and the capacity to evolve to about  
637 17,000 (in the winter) or no such constraint (see Methods). For both constrained and  
638 unconstrained simulations, wheel running for HR and control lines diverge recognizably  
639 at least by generation 6 (Fig 5A), consistent with the selection experiment. The  
640 replicate HR lines for unconstrained and constrained models appear fairly similar for  
641 earlier generations (Figs 5B-C), presumably because mice are not widely achieving  
642 constrained running levels. As expected, the among-line variation for control lines  
643 increases gradually across generations. For the HR lines, among-line variance does  
644 not increase to a noticeable extent and potentially even diminishes by later generations,  
645 a result that is also consistent with the selection experiment (Garland, Jr., Kelly, et al.  
646 2011).

647

648 **Fig 5. Power simulations considering a constraint.** Simulated running levels for (A)  
649 mean running of 4 non-selected control lines (blue) for 200 simulations compared with  
650 mean running of 4 HR lines under 100 unconstrained simulations (dark red) and 100  
651 constrained simulations (light red). (B) Individual HR and control lines single  
652 unconstrained simulation. (C) Individual HR and control lines single constrained  
653 simulation. Black lines show the mean narrow-sense heritability for four lines within  
654 each linetype for (D) 100 control line simulations (arbitrarily the first 50 from each  
655 model), (E) 100 simulations for HR lines with unconstrained running, and (F) 100  
656 simulations for HR lines with constrained running (16,000 revolutions). Standardized

657 selection differentials (calculated within family and sex) from simulations for (G) control  
658 lines, (H) HR lines without a constraint, and (I) HR lines with a constraint. The red line  
659 represents the mean of all heritability and selection differentials (N = 100). Note  
660 different axes for panel G versus H and I.

661

662 The calculated heritability (slope of the regression of offspring [generation 1] on  
663 midparent [generation 0]) for all 200 simulations for control lines indicate that our  
664 parameters resulted in a narrow-sense heritability of about 0.3621 (N=8,000 families).

665 For individual lines, the estimated heritability for successive generations was highly  
666 variable, as would be expected with such small sample sizes (10 families/line).

667 However, the means clearly indicate a slow loss of heritability in the control lines and a  
668 more rapid loss in the HR lines, although values never go to zero (Figs 5E-F),  
669 consistent with the selection experiment (Careau et al. 2013).

670 The standardized selection differentials (calculated within family and sex) for the  
671 unconstrained model remained very consistently around 0 for the control lines (Fig 5G)  
672 and 1.2-1.3 for the HR lines (Figs 5H-I). However, the constrained selection differential  
673 is on average 0.016 below the unconstrained differential. Although slight, this difference  
674 remains consistent across nearly all generations (graph not shown). In the actual  
675 selection experiment, selection differentials declined across generations (Careau et al.  
676 2013).

677 Under both models, Type I error rate for  $\alpha = 0.05$  when comparing allele  
678 frequencies of HR with C lines was deflated at generation 0, regardless of the effect  
679 size for the locus. Type I error ranged from 0.0313 to 0.0420 with no preference for any

680 effect size (S1 Table). This relatively low power when comparing the HR and control  
681 lines (when the line itself is the experimental unit) has been documented previously with  
682 simulations for both genetic data and phenotypes (Hillis et al. 2020; Castro et al. 2021).

683 As expected, power to detect differentiation between the HR and C lines  
684 increased across generations, but never exceeded 0.057 for any generation for either  
685 model. Comparing models at each generation indicates that power is significantly  
686 higher under the unconstrained model by generation 15, although the difference is trivial  
687 (0.0015 with P=0.0127) (Table 4). This differential in power increased through  
688 generation 50, when it reached 0.0055 (P=2.63E-15), before beginning to diminish with  
689 later generations. Although the information in Table 4 does not tell us about the power  
690 to detect loci based on effect size (see numbered list below), it does establish that we  
691 expect more total selection signatures at generation 61 than 22 (see Discussion).

692

693 **Table 4. Simulated statistical power for constrained and unconstrained models**

Generation	Unconstrained power	P-value for Unc.		P-value for C.	
		Comparing present generation to previous generation <sup>a</sup>	Constrained power	Comparing present generation to previous generation <sup>a</sup>	P-value for Constrained vs Unconstrained <sup>b</sup>
0	0.0408 <sup>c</sup>	NA	0.0409 <sup>c</sup>	NA	0.8436
5	0.0417	0.1192	0.0423	0.0365	0.3688
10	0.0429	0.0379	0.0429	0.2685	0.9817
15	0.0447	0.0060	0.0432	0.6693	0.0127
20	0.0456	0.1272	0.0425	0.1891	2.69E-07
22	0.0461	0.4333	0.0425	0.9470	2.31E-07
25	0.0465	0.5290	0.0434	0.1721	1.58E-06
30	0.0484	0.0042	0.0447	0.0374	4.61E-08
35	0.0503	0.0041	0.0463	0.0155	4.51E-09
40	0.0527	0.0006	0.0476	0.0423	5.31E-13
45	0.0544	0.0160	0.0497	0.0009	6.42E-11
50	0.0557	0.0584	0.0502	0.4380	2.63E-15
55	0.0562	0.4050	0.0510	0.2169	2.40E-13
60	0.0565	0.7307	0.0512	0.8162	3.70E-13
61	0.0563	0.8387	0.0512	0.9147	6.27E-13

694  
695 <sup>a</sup>These p-values are calculated using a T-test assuming unequal variance comparing the power  
696 of the generation for that line to the previously listed generation (e.g., unconstrained power at  
697 generation 5 compared to generation 0 has a p-value of 0.1317).  
698 <sup>b</sup>P-value from a T-test (unequal variance) comparing the power of the 100 unconstrained  
699 simulations to the 100 constrained simulations  
700 <sup>c</sup>Type I error rate.

701  
702 The average Pearson correlation between p-values across 2,096 loci for  
703 generation 22 and 61 for the constrained model ( $r = 0.3843$ ) was not statistically  
704 different from that for the unconstrained model ( $r = 0.3889$ : total  $N = 200$ , unpaired-T = -  
705 1.6938,  $P = 0.0919$ ). In the unconstrained model, 33.3% of the loci significantly  
706 differentiated at generation 22 ( $\alpha = 0.05$ ) were still differentiated at generation 61,  
707 versus 31.5% under the constrained model (unpaired-T of percentage of loci  
708 consistently different for generations 22 and 61 at  $p \leq 0.05$  for 200 simulations = -  
709 2.4269,  $P = 0.01614$ ). This consistency of about 1/3 is more than 3 times greater than  
710 for the real data (9.12% for T-tests), which mirrors the drop in the correlation of p-values  
711 between generations 22 and 61 (Fig 2B,  $r = 0.0898$ ). Incorporation of sampling error  
712 into the constrained model lowered the correlation between p-values to 0.2695 and the  
713 proportion of loci significant at generation 22 still significant at generation 61 to 22.7%.  
714 Comparisons of power to detect differentiation between the HR and C lines in  
715 relation to effect size of locus and generation under two simulation models (S2 Table)  
716 indicates:  
717 1. power increased with effect size, as expected;  
718 2. at generation 22, power was greater in the unconstrained model for loci with effect  
719 sizes 25.6 or above;

720 3. at generation 61, power was greater under the unconstrained model for loci with  
721 effect sizes 12.8 and above (excluding the largest effect size, 204.8);  
722 4. under both models, power was consistently greater at generation 61, except for the  
723 largest effect size, where the power is reversed.

724

725 [Comparison of selection signatures in generations 22 and 61 for pooled data](#)

726 When the average allele frequencies of SNPs within regions identified by strict culling at  
727 generation 22 (this study) are compared to the average AF of those loci at generation  
728 61, an increase in among-line variance is apparent for generation 61, within both the HR  
729 and C linetypes (Fig 6A). All else being equal, this increase in among-replicate variance  
730 should lower the statistical power to detect differentiation between the HR and C  
731 linetypes. In agreement with this expectation, most of these strict regions at generation  
732 22 (Table 2) are no longer significantly differentiated at generation 61 (Table 2 and Fig  
733 6A). However, several of the 5 strictly culled regions at generation 61 also show some  
734 evidence of differentiation at generation 22 (Fig 6E). Although strict culling methods  
735 exclude regions identified with generation 61 AF analyses, regions implicated in  
736 generation 22 with  $FDR = 0.01$  culling alone do have considerable overlap with some of  
737 these 5 regions identified at generation 61. Generation 61 regions 3, 4, and 5 (Table 2)  
738 were significant at  $FDR = 0.01$  for all three analyses at generation 22 (S3 Table).

739

740 **Fig 6. Power simulations considering a constraint.** (A) Allele frequencies of regions  
741 identified as significant (via strict culling) at generation 22 (Table 2) (excepting region  
742 18, for which no loci were available in the generation 61 data). Line plot examples are

743 provided for generation 22 (B) region 1, (C) region3, and (D) region 10, with generation  
744 22 on left and generation 61 on right. (E) Allele frequencies of regions identified as  
745 significant at generation 61. (F) Line plot example for generation 61 region 6. SNPs  
746 included have a nominal  $p < 0.05$  at generation 22 and any SNPs at generation 61 which  
747 matched the generation 22 SNPs (shared loci, see Table 1).

748

749 For some of the regions identified as significant at generation 22, differentiation  
750 may have been lost by generation 61 as result of a single line diverging from the others  
751 (for example, line 3 in region 1 or line 7 in region 3 [Figs 6B-C]). In general, mean  
752 differences at generation 22 are much smaller than at 61, but also with much less  
753 among-line variance. A particular example of this includes region 10 (Fig 6D), which is  
754 the only region identified at generation 22 (after strict culling) to continue to be detected  
755 as differentiated at generation 61 (see Discussion).

756

#### 757 [Possible biological function of generation 22 differentiated regions](#)

758 A total of 79 genes (including predicted genes and miRNA) were identified using “strict”  
759 culling of generation 22 regions. These were insufficient for powerful ontology tests and  
760 so regions containing at least five differentiated loci at FDR  $< 0.01$  were included,  
761 bringing the total number of included genes to 345 (S4 Table). Of these 345  
762 differentiated genes, 285 were recognized by the Panther database (used by the Gene  
763 Ontology Resource) and used for identifying potential biological function. Those not  
764 recognized were generally miRNA, predicted, or olfactory genes.

765                   GO Biological Functions implicate antifungal innate immune response, sensory  
766 perception of smell, and embryonic skeletal system morphogenesis. The antifungal  
767 enrichment appears to be the result of a group of C-type lectin genes found on  
768 chromosome 6 (chr6:122,815,876-124,446,843). It may not be a coincidence that this  
769 region also includes a group of vomeronasal genes contributing to the sensory  
770 perception of smell term. The genes that implicate the embryonic skeletal  
771 morphogenesis term include a cluster of *Hoxb* genes found on chromosome 11  
772 (chr11:93,129,916-96,570,699).

773                   A few genes specifically from the 79 found in the “strict” regions merit mention for  
774 their relevance to the running phenotype, including *Cited2* (adrenal cortex formation),  
775 *Rbm24* (positive regulation of skeletal muscle fiber differentiation), and *Dspp* (negative  
776 regulation of bone development).

777

778

779

780

781

782 [Discussion](#)

783 [Overview](#)

784 Previously, whole-genome sequence data for individual mice at generation 61 of the

785 High-Runner mouse selection experiment were used to identify genomic regions

786 differentiated between HR and control lines. Thirteen of these were termed "consistent"

787 because they appeared with three different analytical methods (Hillis et al. 2020).

788 These 13 regions contained genes associated with known phenotypic differences

789 between the HR and control lines and intuitive associations with running ability and/or

790 motivation/reward systems. However, given that the HR lines had begun to reach

791 selection limits around generations 17-27, depending on line and sex (Careau et al.

792 2013), tens of additional generations, with continuing random genetic drift, could have

793 obscured many selection signatures. Therefore, in the present study, we analyzed

794 allele frequencies for the lines sampled at generation 22, based on DNA pooled by line.

795 These analyses of generation 22 identify hundreds of genomic regions differentiated

796 between the HR and C lines (FDR = 0.01), despite using pooled sequence data rather

797 than sequences for individual mice (Xu and Garland 2017). We then reanalyzed the

798 data from generation 61 as allele frequencies by line, to mimic the data available for

799 generation 22, and found that the regions identified as differentiated at generation 61

800 are, at best, weakly differentiated at generation 22. Nevertheless, both generations'

801 differentiated regions contain genes that make biological sense for wheel-running

802 behavior. Below, we discuss (1) implications of the differences in data type between

803 generations 22 and 61, (2) possible statistical and biological explanations for the

804 differences in identified regions, and (3) genes and biological systems highlighted by the  
805 genomic regions identified by generation 22 analyses (after strict culling).

806

807 **Differences in selection signatures at generations 22 and 61**

808 We expected estimates of selection signatures to be similar at generations 22 and 61,  
809 based on the fact that the HR lines had mostly reached selection limits by generation 22  
810 (Careau et al. 2013), such that the most biologically important loci would have gone to  
811 fixation or at least reached equilibria across most or all of the HR lines. In agreement  
812 with this expectation, of the 13 "consistent" regions identified by Hillis et al. (2020) for  
813 generation 61 (using individual mouse data), 8 were still identified by at least one of the  
814 tests (FDR = 0.01) using the generation 61 genotypes pooled into allele frequencies per  
815 line. Generation 22 analyses of pooled sequence data identified 7 of the 13 consistent  
816 regions (although several of these regions were only detected by a few SNPs: S3  
817 Table). Interestingly, the consistent region on chromosome 14 was more strongly  
818 detected at generation 22 than at generation 61 using pooled sequence analyses (Table  
819 2).

820 On the other hand, the strongest selection signatures observed at generation 61  
821 with the data treated as pooled sequences are not among the strongest ones observed  
822 at generation 22 (based on number of SNPs detected and their p-values), despite  
823 continued selection on the HR lines. For example, region 1 of the generation 61 "strict"  
824 culling pooled analyses (chr3:51,199,110-51,602,693) included 124 SNPs (Table 2),  
825 and whereas generation 22 analyses did not detect any of these loci as significant.  
826 Another example, region 3 of the generation 61 pooled analyses (chr9:41,413,436-

827 42,478,817) included 1,277 SNPs (FDR = 0.01), but none of the 647 shared loci were  
828 identified by in the generation 22 analyses (S3 Table). When directly comparing SNPs  
829 differentiated at FDR = 0.01, we see only a single SNP of overlap for the WRT test.

830 In addition to the differences in individual SNP results, a 17-fold greater number  
831 of regions was identified by generation 22 analyses than generation 61 pooled analyses  
832 at FDR = 0.01 (Table 1). This ratio applies to all statistical tests and the complete SNP  
833 analyses for each generation, as well as the analyses of SNPs shared by the two  
834 generations. Moreover, the SNPs identified at generation 61 were clustered into far  
835 fewer regions (Table 1). Broadly, this difference in numbers of selection signatures  
836 have at least two possible explanations, which are not mutually exclusive: (1)  
837 differences in data type, quality, and quantity; (2) biological differences between  
838 generations 22 and 61.

839

840 1) Differences in data type, quality, quantity, and sampling error

841 Our power to detect differentiation in allele frequencies should have been lower for  
842 generation 22 than for generation 61 (Figs 4C and D). As also noted in the Methods,  
843 the estimates for SNP allele frequencies per line at generation 61 were based on ~10  
844 mice/line sampled and an 12X average read depth per mouse, yielding a total of  
845 5,932,148 variable SNP loci (Hillis et al. 2020). For generation 22, pooled sequencing  
846 was done with ~20 mice/line and an average read depth of 24X, yielding 4,446,523  
847 variable SNPs (Table 1). Generally, with an average read depth of 12X per mouse,  
848 both alleles will be represented for each mouse (i.e., 20 alleles per line) for generation  
849 61 allele frequencies. However, with 24X average read depth for generation 22,

850 simulations involving sampling alleles with replacement show that generation 22 is  
851 prone to vary more from the actual population allele frequency (Figs 4A and B). Thus,  
852 the much greater number of differentiated SNPs and chromosomal regions detected at  
853 generation 22 would not appear to be simply a function of greater statistical power  
854 versus generation 61. Thus, we now consider possible biological explanations.

855

856 2) Biological differences

857 One way to highlight the differences in selection signatures detected at generations 22  
858 and 61 is to note that of the differentiated regions detected for generation 61, two of  
859 them contain hundreds of statistically significant SNPs (FDR = 0.01) shared between  
860 the generation 22 and 61 data sets. Despite this, those two regions are not among the  
861 more differentiated regions in the Manhattan plots (Fig 3A and B).

862 What biological explanations might account for such discrepancies? One  
863 possibility is a physiological constraint that eliminates the need for all loci favorable to  
864 wheel running to be maintained at high frequencies once a selection limit is reached  
865 (see verbal model in Hillis and Garland, Jr. 2022). We consider this from the  
866 perspective that many complex traits are influenced by hundreds or thousands of loci  
867 (Wood et al. 2014; Long et al. 2015). Voluntary exercise behaviors would likely be  
868 among them, given that they incorporate numerous physiological and morphological  
869 traits related to ability (e.g., cardiac muscle, skeletal muscle, bone, metabolism, water  
870 and temperature homeostasis) as well as aspects of motivation and reward (e.g.,  
871 dopamine signaling, chemosensory systems) (Lightfoot et al. 2018; Wang et al. 2022).

872           Although biological constraints can be defined in various ways (Garland et al.  
873   2022), in the present context, a constraint would be anything that limits the maximum  
874   revolutions that an individual mouse can run during the testing period. Previously, we  
875   discussed how different unique responses to identical selection criteria (i.e., “multiple  
876   solutions”) could occur and referenced constraints as a potential explanation (Hillis and  
877   Garland, Jr. 2022). To utilize and expand on their example, suppose that mice are  
878   subject to a constraint on wheel running caused by joint pain: they stop running when  
879   the pain becomes intolerable. In this scenario, joint pain is sufficient to limit wheel  
880   running and it serves as a “weak link” or single limiting factor in the biological systems  
881   required for high wheel running. Then suppose 10 alleles located at 10 independent  
882   biallelic loci, with entirely additive effects, are capable of increasing wheel running.  
883   Suppose further that only five such alleles are needed to achieve the maximum amount  
884   of wheel running permitted by joint pain. Under this model, if selection acts on a  
885   population to increase running, then (1) fixation of the favorable allele at any five of the  
886   10 loci will coincide with a selection limit determined by pain tolerance, (2) none of the  
887   alleles at any of the 10 loci must be fixed to reach the pain-determined limit, (3) more  
888   than 5 favorable alleles could be maintained at intermediate allele frequencies, and (4)  
889   as long as enough favorable alleles are maintained for the selection limit, some  
890   favorable alleles can be lost without detriment to wheel running. These factors allow for  
891   substantial variation among the replicate lines and considerable flexibility for change  
892   within a given line, even for favorable wheel-running alleles at the selection limit. This  
893   possibility of “genetic churn” beyond a selection limit that is caused by a physiological  
894   constraint also implies that genotype-to-phenotype maps (Travisano and Shaw 2013;

895 Zamer and Scheiner 2014; Porto et al. 2016; Zinski et al. 2021) may be moving targets  
896 and hence difficult to identify. Therefore, we used simulations to compare power to  
897 detect and consistency in detected selection signatures, both with and without a  
898 physiological constraint.

899

#### 900 [Allele frequency divergence over time](#)

901 The relatively large number of SNPs identified to be divergent between generations 22  
902 and 61 illustrate a shift between the generations. Furthermore, regions with the  
903 greatest divergence between the generations align closely with some of the previously  
904 identified regions, particularly those identified by the generation 61 mixed model  
905 analyses (Table 3). By generation 22, the lines had not had as much time to evolve as  
906 much separation in allele frequencies between HR and C as by generation 61. As a  
907 result, the differences between the linetypes are generally between -0.5 and 0.5 at  
908 generation 22 (Fig 3C), whereas this difference expands to between -1 and 1 for  
909 generation 61 (Fig 3D). Given that the significant regions identified at gen 61 are  
910 typically going to be those whose difference is near -1 or 1, most likely those same  
911 regions at gen 22 had differences within the -0.5 and 0.5 range. Therefore, a growing  
912 difference in HR and C allele frequency had to have occurred, in our data, between gen  
913 22 and gen 61 to observe a significant difference at 61. Such loci would naturally be  
914 among the more significant in the difference of t-tests over time.

915

916 Simulations comparing power under constrained versus unconstrained models  
917 These simulations were conducted to test whether a physiological constraint on the  
918 phenotype of wheel-running behavior could reduce the consistency of loci identified at  
919 different generations of a selection experiment. To better simulate realistic phenotypic  
920 variance within the population, both the wheel-running and constraint phenotype  
921 simulations were based on equations with both genetic and environmental sources of  
922 variance, such that both could evolve.

923  
924 Similarities between the constrained model and real data  
925 The constrained model appears to better reflect what we observe in the actual response  
926 to selection. This is due to the lack of a clear selection limit achieved under the  
927 unconstrained model (Fig 5A). Although the response to selection diminishes over time  
928 (likely due to the reduction in heritability: Fig 5E), a clear plateau is not apparent. In  
929 contrast -- as must be the case -- a clear plateau occurs under the constrained model.

930  
931 Correlation between generations 22 and 61 p-values  
932 For the tests comparing allele frequencies at each of 2,096 loci between the HR and C  
933 lines, the correlation of p-values between generations 22 and 61 was statistically lower  
934 in the constrained model ( $r=0.3843$ ) as compared to the unconstrained model  
935 ( $r=0.3889$ ), though still 4x higher than for the real genomic data ( $r=0.0909$ ) (Fig 2).  
936 Even with the inclusion of sampling error, the correlation ( $r=0.2695$ ) is nearly 3x greater  
937 than for the real genomic data, which indicates that other factors (e.g., gene  
938 interactions) must be contributing to the differences between the generations.

939           In spite of the similarity in the correlation of p-values between generations for  
940    simulations, the between-generation consistency of detected selection signatures was  
941    slightly but statistically greater under the unconstrained (33.3%) than under the  
942    constrained model (31.5%). This difference may be due to the constrained model  
943    having very slightly (~0.016) though consistently lower selection differentials (Figs 5G  
944    versus 5H), which could lead to less fixation of favored alleles. However, the relatively  
945    small difference between models in consistency of selection signatures is not enough to  
946    explain the large differences in the real data between generations 22 and 61 (Table 2).  
947    The inclusion of sampling error into the estimates decreased the 31.5% consistency  
948    between generations 22 and 61 differentiated loci to 22.7%. This level of consistency  
949    with simulated data remains more than 2-fold higher than for the real data (9.12%), thus  
950    implicating the presence of additional factors that reduce consistency in the real data  
951    (e.g., epistatic effects).

952           Overall, our simulations fail to demonstrate why we observe a 17X drop in  
953    significant regions from generation 22 to 61 (Table 1), implying instead that we should  
954    detect more at generation 61 than at 22 (Table 4).

955

#### 956    Effect sizes of loci

957    Under both models, more loci were detected at later generations (Table 4). However,  
958    the power to detect loci with the largest effect size was much higher at earlier than later  
959    generations (S2 Table and S1 Fig). This pattern makes sense in consideration of the  
960    factors that affect the average difference in allele frequency between the HR and control  
961    lines and the variance among replicate lines within linetypes. Drift will generally

962 increase the variance among lines with each generation. The allele frequencies in the  
963 simulated control lines will be affected only by this drift. Allele frequencies in the HR  
964 lines will be affected both by drift and selection, where selection will have stronger  
965 effects at loci with larger effect sizes. This results in something of a race between  
966 selection increasing the difference in allele frequencies between HR and control lines,  
967 while drift increases among-line variance for both HR and control lines. For loci with  
968 small effect sizes, drift will have a relatively greater influence over allele frequencies  
969 than selection at any generation, and thus detection rates never vary far from the Type I  
970 error rate, i.e., power is virtually zero (S2 Table). Loci with large effect sizes, however,  
971 are able to differentiate rapidly, often leading to fixation of the favored allele in our  
972 simulations (S1 Fig). Even after fixation in the HR lines, drift is still able to increase  
973 allele-frequency variance among the control lines (potentially to the point of fixing loci  
974 for opposite alleles), thus further reducing the power to detect any differentiation. Thus,  
975 the power to detect signatures of selection should increase the most rapidly across  
976 generations for loci with the largest effect sizes, but power is also expected to decline  
977 after fixation of the favored alleles in the HR lines and with continuing increase in  
978 variance among the control lines (S1 Fig).

979 That the power to detect a locus as differentiated is correlated with its effect size  
980 is unsurprising. For example, under the unconstrained model the power to detect  
981 selection signatures for loci with 0.4 effect size is about 16.6-fold less than the power to  
982 detect loci with 204.8 effect size and 12.8-fold less for the constrained model  
983 (generation 22). This gap diminishes to about 7.7-fold difference (both models) by  
984 generation 61 (S2 Table), presumably due to the reasons described in the previous

985 paragraph. However, the 0.4 effect size loci are far more numerous than the 204.8  
986 effect size loci ( $N = 720$  and 8, respectively). Consequently, the number of 0.4 effect  
987 size loci detected as significant is nearly 5-fold greater (unconstrained) and more than  
988 7-fold greater (constrained) than the number of 204.8 effect size loci detected. The  
989 most notable difference between the constrained and the unconstrained models is that  
990 at generation 22 the unconstrained model yielded substantially more power than the  
991 constrained model for loci with the largest effect sizes (0.724 to 0.485, respectively;  
992 unpaired t-test,  $P=4.01E-22$ ). This would imply that constraints may have a substantial  
993 impact on the ability to detect selection at loci with the greatest effect sizes, a result that  
994 deserves further study.

995 For identifying possible biological functions, we would ideally focus on loci with  
996 relatively large effect size, as these will have the most direct influence on the phenotype  
997 and may serve as potential targets for future functional studies. We have no information  
998 on effect sizes of SNPs or regions detected as differentiated for our real data. The  
999 relative proportions of low- and high-effect size loci among the detected selection  
1000 signatures in the real data will likely vary from our simulations, depending on the actual  
1001 distribution of those effect sizes and other factors. However, the simulations do suggest  
1002 that we may have numerous small-effect size loci among our detected selection  
1003 signatures. The inclusion of the “strict” culling method was meant to prioritize regions  
1004 that would have large effect sizes. Having more loci that are differentiated and linked  
1005 together would be expected from those regions under strong selection because  
1006 recombination would have fewer generations to break up linked base pairs before the  
1007 region becoming fixed in the HR lines. As the simulations have so many more loci with

1008 small effect sizes, at generation 0, when we compare the lowest p-value produced for  
1009 each simulation for the 0.4 effect size we tend to see lower p-values than loci with 204.8  
1010 effect size simply because of more opportunities to produce a low p-value. However,  
1011 generation 22 appears to be better for detecting a greater proportion of selection  
1012 signatures from loci with large effect sizes as the relative proportion on large effect size  
1013 loci appears to be higher (S2 Table).

1014

1015 [Possible biological functions of generation 22 differentiated regions](#)  
1016 Ontology analyses identified biological processes that can be grouped into three  
1017 categories: sensory perception of smell, antifungal innate immune response, and  
1018 embryonic skeletal system morphogenesis. Of these, the system that is most  
1019 consistent between generations 22 and 61 is the perception of smell, which was among  
1020 the mostly clearly differentiated systems at generation 61 (Hillis et al. 2020). As was  
1021 discussed by Hillis et al. (2020), the experimental procedure for measuring wheel  
1022 running, for logistical reasons, involved mice being placed on wheel over three batches  
1023 and mice in batches 2 and 3 are placed on wheels which still smell of the previous  
1024 mouse (Swallow, Garland, Jr., et al. 1998). Evidence of an evolutionary response to  
1025 this is visible in the HR lines in that HR mice will run at very different speeds if on a  
1026 wheel that is clean, previously traversed by a male, or previously traversed by a female  
1027 (Dewan et al. 2019). Alterations in the transcriptome also indicates changes in olfactory  
1028 and vomeronasal systems (Nguyen et al. 2020). Taken together, these results indicate  
1029 that perception of smell may be a notable factor in their motivation for running on the  
1030 wheels and also consistent with the idea that motivation is expected to evolve before

1031 ability (Garland, Jr. et al. 2016; Khan et al. 2024). Interestingly, although both  
1032 generations demonstrate evolution in genomic regions association with olfaction and  
1033 vomeronasal, the regions implicated in each generation are different, with an exception  
1034 of the region on chromosome 14 (chr14:52,115,206-53,776,455), which was identified  
1035 by the generation 22 WRT analyses and the generation 61 mixed model analyses  
1036 (Table 2). However, additional studies should be done to address the effects of  
1037 olfactory/vomeronasal systems more directly on running behavior of the HR mice. This  
1038 could be done with ablation procedures on HR and C mice and observing changes in  
1039 running behavior. The antifungal ontology term is possibly a hitchhiker with the  
1040 vomeronasal genes also present in the differentiated region (chr6:122,815,876-  
1041 124,446,843).

1042         Ontology analyses also indicated embryonic skeletal system development as  
1043 result of a group of *Hoxb* genes within a differentiated region. If these *Hoxb* are the  
1044 driving force underlying the many skeletal differences that have been documented  
1045 between HR and C lines (Garland, Jr. and Freeman 2005; Kelly et al. 2006; Middleton  
1046 et al. 2008; Middleton et al. 2010; Wallace et al. 2010; Wallace et al. 2012; Castro and  
1047 Garland, Jr. 2018; Copes et al. 2018; Schwartz et al. 2018), then this is an exciting  
1048 discovery because it would represent a response to selection in a group of genes  
1049 known to be evolutionarily influential in body patterning and development (Stratford et  
1050 al. 1999). However, whether the *Hoxb* genes are the cause of skeletal differentiation is  
1051 unclear. Although *Hoxb* genes may play a role in these changes, they are far from the  
1052 only candidates. GO term “skeletal system development” includes 7 additional non-*Hox*  
1053 genes, including *Phospho1*, *Col1a1*, and *Mbtd1*, which are all located in the same

1054 differentiated region as the *Hox* genes. Furthermore, individual loci demonstrating the  
1055 greatest differentiation do not appear to be in *Hox* genes themselves or their regulatory  
1056 regions. Even if *Hox* genes are a hitchhiker in a region with other genes more directly  
1057 targeted by selection due to their skeletal effects, exploring potential side effects of this  
1058 evolution would be of interest. Expression analyses during developmental stages when  
1059 these genes are most active may provide insight into how *Hox* genes may be altered in  
1060 the HR mice.

1061

1062 Other genes of potential interest

1063 The 79 genes included in top regions also contain a few of particular note:  
1064 *Cited2*, *Rbm24*, and *Dspp*. Each of these genes is associated with ontologies and  
1065 phenotypes that have been identified as differentiated between the HR and C mice.  
1066 *Cited2* is a gene whose knockout (KO) has been associated with alterations in brain and  
1067 heart morphology (Barbera et al. 2002; Bamforth et al. 2004; MacDonald et al. 2008)  
1068 and has also been associated with adrenal development (Val et al. 2007). As noted in  
1069 the introduction, HR mice have larger brains and hearts than C mice (Kolb et al. 2010;  
1070 Kolb, Rezende, et al. 2013; Kolb, Kelly, et al. 2013; Kelly et al. 2017). Additionally,  
1071 adrenal corticosterone levels were found to be different between the linotypes (Malisch  
1072 et al. 2007; Garland, Jr. et al. 2016). *Rbm24* is a gene associated with skeletal muscle  
1073 fiber differentiation, particularly during regeneration following injury (Cardinali et al.  
1074 2016; Zhang et al. 2020; Grifone et al. 2021). The HR and C lines have demonstrated  
1075 differences in muscle fiber types within muscles important for wheel running such as the  
1076 gastrocnemius (Syme et al. 2005; Guderley et al. 2008; Castro et al. 2022). However,

1077 differential response to muscle injury has not been found between the linotypes (Kay et  
1078 al. 2022). Lastly, *Dspp* was identified among the differentiated genes. This gene has  
1079 been associated with development of long bones (such as femurs) and cortical and  
1080 trabecular bone thickness (Verdelis et al. 2008; Jani et al. 2016). The HR and C mice  
1081 have shown various differences in bone morphology (see Introduction).

1082

### 1083 Limitations and conclusions

1084 Some of the limitations of the present study include trying to compare results of pooled  
1085 genome sequencing (generation 22) to individual mouse sequencing (generation 61:  
1086 Hillis et al. 2020). Though the alleles of the individual mice can be combined to imitate  
1087 pooled genome sequences, the differences in number of mice sampled and sampling  
1088 error make comparisons problematic (see Methods). This is illustrated by the decrease  
1089 in p-value correlations (between generations 22 and 61) as compared to both the  
1090 unconstrained and constrained simulations. Nevertheless, as argued above, neither the  
1091 increase in number of regions detected as differentiated at generation 22 nor the lack of  
1092 correspondence between detected regions at generations 22 and 61 can be explained  
1093 solely by methodological differences.

1094 The constraint simulations have their own limitations in that they do not account  
1095 for male vs female running differences (females run more than males)(Careau et al.  
1096 2013). In addition, dominance, epistasis, and gene-environment interactions were not  
1097 considered. Exclusion of these features may be why we were unable to achieve  
1098 realistic levels of among-line variation, particularly among the High Runner lines. This  
1099 model also does not include linkage disequilibrium or realistic rates of recombination.

1100 Additionally we do not include reduction in breeding success across generations, which  
1101 may explain the drop in selection differential observed by Careau et al. (2013). Lastly,  
1102 we did not explore the potential effects of relaxing selection for four generations, as  
1103 when the mice were moved from Wisconsin to California (see Introduction). A cluster of  
1104 generations of no selection in the HR lines could allow for some drift of the favored  
1105 alleles.

1106 Although, we are unsure as to why we see so many regions at FDR = 0.01 that  
1107 do not correspond to the generation 61 findings by Hillis et al. (2020), our simulations  
1108 suggest that regions with the strongest effect sizes on wheel running are likely to be  
1109 among the generation 22 regions. Given the statistical significance and number of  
1110 SNPs identified in our “strictly” culled differentiated regions, these regions are most  
1111 likely to have had the greatest impact on wheel running at the start of the selection  
1112 experiment. Among these regions are genes related to olfactory/vomeronasal function,  
1113 reward pathways, and a miRNA cluster that has been associated with energy  
1114 homeostasis in neonatal development. All of these associations make sense based on  
1115 known phenotypic differences between the HR and control lines (see Introduction).

1116 Future directions might include more complex simulations (e.g., see Baldwin-  
1117 Brown et al. 2014; Stephan 2016; Castro et al. 2019), which may better help to explain  
1118 the 17X increase in regions detected at generation 22. Including genomic data from  
1119 more generations (especially from the base population, generations near to but before  
1120 the selection limit, and current generations [i.e., around 100]) may provide more clarity  
1121 regarding how the response to selection changes across phases of the selection  
1122 response (cf. Rose et al. 2005; Castro et al. 2021). Analyses using all loci and a kinship

1123 matrix would enable determination of some interactions between genes. Functional  
1124 analyses, such as knockouts of some of the genes whose alleles appear to have been  
1125 favored by selection, may provide direct evidence of influence on wheel-running  
1126 behavior (e.g., Schmidt et al. 2008; Chaouloff et al. 2011; MacKay et al. 2019).  
1127 Furthermore, analyses of other physiological aspects of these KO mice may help to  
1128 better understand the mechanisms by which the gene influences wheel running.

1129

1130

1131 [Acknowledgements](#)

1132 We would like to thank Dr. Lei Yu for help with SNP calling, Layla Hiramatsu for  
1133 assistance in sample collection, and Dr. Zhenyu (Arthur) Jia, Dr. Shizhong Xu, Dr. Frank  
1134 Chan, and Dr. Tony Long for comments and suggestions for this study.

1135

1136 [Competing of interests](#)

1137 The authors have no competing interests.

1138

1139 [Author Contributions](#)

1140 Conceptualization, D.A.H., Lir.Y., F.P.M.dEV., D.P., T.G.; investigation, D.A.H., Lir.Y.,  
1141 G.M.W., F.P.M.dEV., D.P., T.G.; software, D.A.H., Lei.Y.; formal analysis, D.A.H., Lir.Y.,  
1142 T.G.; writing – original draft, D.A.H., T.G.; writing review and editing, all authors.

1143

1144 [Data Availability](#)

1145 Generation 61 data were made available by Hillis et al. (2020) and can be found at  
1146 <https://doi.org/10.25386/genetics.12436649>. Generation 22 fastq files are available on  
1147 the SRA database, accession = PRJNA758905  
1148 (<https://www.ncbi.nlm.nih.gov/sra/?term=PRJNA758905>).

1149 [References](#)

1150 Ahrens WH, Cox DJ, Budhwar G. 1990. Use of the arcsine and square root  
1151 transformations for subjectively determined percentage data. *Weed Sci.* 38:452–  
1152 458.

1153 Al-Murrani WK, Roberts RC. 1974. Genetic variation in a line of mice selected to its limit  
1154 for high body weight. *Anim. Sci.* 19:273–289.

1155 Baldi P, Long AD. 2001. A Bayesian framework for the analysis of microarray  
1156 expression data: regularized t-test and statistical inferences of gene changes.  
1157 *Bioinformatics* 17:509–519.

1158 Baldwin-Brown JG, Long AD, Thornton KR. 2014. The power to detect quantitative trait  
1159 loci using resequenced, experimentally evolved populations of diploid, sexual  
1160 organisms. *Mol. Biol. Evol.* 31:1040–1055.

1161 Bamforth SD, Bragaña J, Farthing CR, Schneider JE, Broadbent C, Michell AC, Clarke  
1162 K, Neubauer S, Norris D, Brown NA, et al. 2004. Cited2 controls left-right  
1163 patterning and heart development through a Nodal-Pitx2c pathway. *Nat. Genet.*  
1164 36:1189–1196.

1165 Barbera JPM, Rodrigues TA, Greene NDE, Weninger WJ, Simeone A, Copp AJ,  
1166 Beddington RSP, Dunwoodie S. 2002. Folic acid prevents exencephaly in Cited2  
1167 deficient mice. *Hum. Mol. Genet.* 11:283–293.

1168 Barton NH, Turelli M. 1989. Evolutionary quantitative genetics: how little do we know?  
1169 *Annu. Rev. Genet.* 23:337–370.

1170 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and  
1171 powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B* 57:289–300.

1172 Brown WP, Bell AE. 1961. Genetic analysis of a “plateaued” population of *Drosophila*  
1173 *melanogasteri*. *Genetics* 46:407–425.

1174 Bult A, Lynch CB. 2000. Breaking through artificial selection limits of an adaptive  
1175 behavior in mice and the consequences for correlated responses. *Behav. Genet.*  
1176 30:193–206.

1177 Burke MK, Dunham JP, Shahrestani P, Thornton KR, Rose MR, Long AD. 2010.  
1178 Genome-wide analysis of a long-term evolution experiment with *Drosophila*.  
1179 *Nature* 467:587–590.

1180 Cadney MD, Hiramatsu L, Thompson Z, Zhao M, Kay JC, Singleton JM, Albuquerque  
1181 RL de, Schmill MP, Saltzman W, Garland, Jr. T. 2021. Effects of early-life  
1182 exposure to Western diet and voluntary exercise on adult activity levels, exercise

1183 physiology, and associated traits in selectively bred High Runner mice. *Physiol.*  
1184 *Behav.* 234:113389.

1185 Cardinali B, Cappella M, Provenzano C, Garcia-Manteiga JM, Lazarevic D, Cittaro D,  
1186 Martelli F, Falcone G. 2016. MicroRNA-222 regulates muscle alternative splicing  
1187 through Rbm24 during differentiation of skeletal muscle cells. *Cell Death Dis.*  
1188 7:e2086–e2086.

1189 Careau V, Wolak ME, Carter PA, Garland, Jr. T. 2013. Limits to behavioral evolution:  
1190 the quantitative genetics of a complex trait under directional selection. *Evolution*  
1191 67:3102–3119.

1192 Castro AA, Garland, Jr. T. 2018. Evolution of hindlimb bone dimensions and muscle  
1193 masses in house mice selectively bred for high voluntary wheel-running behavior.  
1194 *J. Morphol.* 279:766–779.

1195 Castro AA, Garland, Jr. T, Ahmed S, Holt NC. 2022. Trade-offs in muscle physiology in  
1196 selectively bred High Runner mice. *J. Exp. Biol.*:jeb.244083.

1197 Castro AA, Rabitoy H, Claghorn GC, Garland, Jr. T. 2021. Rapid and longer-term  
1198 effects of selective breeding for voluntary exercise behavior on skeletal  
1199 morphology in house mice. *J. Anat.* 238:720–742.

1200 Castro JP, Yancoskie MN, Marchini M, Belohlavy S, Hiramatsu L, Kucka M, Beluch NH,  
1201 Rolian C, Chan YF. 2019. An integrative genomic analysis of the Longshanks  
1202 selection experiment for longer limbs in mice. *eLife* 8:e40214.

1203 Chaouloff F, Dubreucq S, Bellocchio L, Marsicano G. 2011. Endocannabinoids and  
1204 motor behavior: CB1 receptors also control running activity. *Physiology* 26:76–  
1205 77.

1206 Copes LE, Schutz H, Dlugosz EM, Judex S, Garland, Jr. T. 2018. Locomotor activity,  
1207 growth hormones, and systemic robusticity: An investigation of cranial vault  
1208 thickness in mouse lines bred for high endurance running. *Am. J. Phys.*  
1209 *Anthropol.* 166:442–458.

1210 Dewan I, Garland, Jr. T, Hiramatsu L, Careau V. 2019. I smell a mouse: indirect genetic  
1211 effects on voluntary wheel-running distance, duration and speed. *Behav. Genet.*  
1212 49:49–59.

1213 Didion JP, Morgan AP, Yadgary L, Bell TA, McMullan RC, Ortiz de Solorzano L, Britton-  
1214 Davidian J, Bult CJ, Campbell KJ, Castiglia R, et al. 2016. *R2d2* drives selfish  
1215 sweeps in the house mouse. *Mol. Biol. Evol.* 33:1381–1395.

1216 Dlugosz EM, Schutz H, Meek TH, Acosta W, Downs CJ, Platzer EG, Chappell MA,  
1217 Garland, Jr. T. 2013. Immune response to a *Trichinella spiralis* infection in house  
1218 mice from lines selectively bred for high voluntary wheel running. *J. Exp. Biol.*  
1219 216:4212–4221.

1220 Dobzhansky T, Spassky B. 1969. Artificial and natural selection for two behavioral traits  
1221 in *Drosophila pseudoobscura*. *Proc. Natl. Acad. Sci. U. S. A.* 62:75–80.

1222 Douhard F, Douhard M, Gilbert H, Monget P, Gaillard J, Lemaître J. 2021. How much  
1223 energetic trade-offs limit selection? Insights from livestock and related laboratory  
1224 model species. *Evol. Appl.* 14:2726–2749.

1225 Dumke CL, Rhodes JS, Garland, Jr. T, Maslowski E, Swallow JG, Wetter AC, Cartee  
1226 GD. 2001. Genetic selection of mice for high voluntary wheel running: effect on  
1227 skeletal muscle glucose uptake. *J. Appl. Physiol.* 91:1289–1297.

1228 Falconer DS. 1989. Introduction to quantitative genetics. 3rd ed. Burnt Mill, Harlow,  
1229 Essex, England : New York: Longman, Scientific & Technical ; Wiley

1230 Garland, Jr. T, Freeman PW. 2005. Selective breeding for high endurance running  
1231 increases hindlimb symmetry. *Evolution* 59:1851–1854.

1232 Garland, Jr. T, Kelly SA, Malisch JL, Kolb EM, Hannon RM, Keeney BK, Van Cleave  
1233 SL, Middleton KM. 2011. How to run far: multiple solutions and sex-specific  
1234 responses to selective breeding for high voluntary activity levels. *Proc. R. Soc. B  
1235 Biol. Sci.* 278:574–581.

1236 Garland, Jr. T, Schutz H, Chappell MA, Keeney BK, Meek TH, Copes LE, Acosta W,  
1237 Drenowitz C, Maciel RC, van Dijk G, et al. 2011. The biological control of  
1238 voluntary exercise, spontaneous physical activity and daily energy expenditure in  
1239 relation to obesity: human and rodent perspectives. *J. Exp. Biol.* 214:206–229.

1240 Garland, Jr. T, Zhao M, Saltzman W. 2016. Hormones and the evolution of complex  
1241 traits: insights from artificial selection on behavior. *Integr. Comp. Biol.* 56:207–  
1242 224.

1243 Garland T, Downs CJ, Ives AR. 2022. Trade-offs (and constraints) in organismal  
1244 biology. *Physiol. Biochem. Zool.* 95:82–112.

1245 Grifone R, Saquet A, Desgres M, Sangiorgi C, Gargano C, Li Z, Coletti D, Shi D-L.  
1246 2021. Rbm24 displays dynamic functions required for myogenic differentiation  
1247 during muscle regeneration. *Sci. Rep.* 11:9423.

1248 Guderley H, Joanisse DR, Mokas S, Bilodeau GM, Garland, Jr. T. 2008. Altered fibre  
1249 types in gastrocnemius muscle of high wheel-running selected mice with mini-  
1250 muscle phenotypes. *Comp. Biochem. Physiol. B Biochem. Mol. Biol.* 149:490–  
1251 500.

1252 Hillis DA, Garland, Jr. T. 2022. Multiple solutions at the genomic level in response to  
1253 selective breeding for high locomotor activity. *Genetics*.

1254 Hillis DA, Yadgary L, Weinstock GM, Pardo-Manuel de Villena F, Pomp D, Fowler AS,  
1255 Xu S, Chan F, Garland, Jr. T. 2020. Genetic basis of aerobically supported

1256 voluntary exercise: results from a selection experiment with house mice.  
1257 *Genetics* 216:781–804.

1258 Hiramatsu L, Kay JC, Thompson Z, Singleton JM, Claghorn GC, Albuquerque RL, Ho  
1259 Brittany, Ho Brett, Sanchez G, Garland, Jr. T. 2017. Maternal exposure to  
1260 Western diet affects adult body composition and voluntary wheel running in a  
1261 genotype-specific manner in mice. *Physiol. Behav.* 179:235–245.

1262 Jani PH, Gibson MP, Liu C, Zhang H, Wang X, Lu Y, Qin C. 2016. Transgenic  
1263 expression of Dspp partially rescued the long bone defects of Dmp1-null mice.  
1264 *Matrix Biol.* 52–54:95–112.

1265 Kay JC, Colbath J, Talmadge RJ, Garland, Jr. T. 2022. Mice from lines selectively bred  
1266 for voluntary exercise are not more resistant to muscle injury caused by either  
1267 contusion or wheel running. Kubis H-P, editor. *PLOS ONE* 17:e0278186.

1268 Kelly SA, Czech PP, Wight JT, Blank KM, Garland, Jr. T. 2006. Experimental evolution  
1269 and phenotypic plasticity of hindlimb bones in high-activity house mice. *J.  
1270 Morphol.* 267:360–374.

1271 Kelly SA, Gomes FR, Kolb EM, Malisch JL, Garland, Jr. T. 2017. Effects of activity,  
1272 genetic selection and their interaction on muscle metabolic capacities and organ  
1273 masses in mice. *J. Exp. Biol.* 220:1038–1047.

1274 Khan RH, Rhodes JS, Girard IA, Schwartz NE, Garland, Jr. T. 2024. Does behavior  
1275 evolve first? Correlated responses to artificial selection for voluntary wheel-  
1276 running behavior in house mice. *Ecol. Evol. Physiol.* 97.

1277 Kolb EM, Kelly SA, Garland, Jr. T. 2013. Mice from lines selectively bred for high  
1278 voluntary wheel running exhibit lower blood pressure during withdrawal from  
1279 wheel access. *Physiol. Behav.* 112–113:49–55.

1280 Kolb EM, Kelly SA, Middleton KM, Sermsakdi LS, Chappell MA, Garland, Jr. T. 2010.  
1281 Erythropoietin elevates  $VO_2$ ,max but not voluntary wheel running in mice. *J. Exp.  
1282 Biol.* 213:510–519.

1283 Kolb EM, Rezende EL, Holness L, Radtke A, Lee SK, Obenous A, Garland, Jr. T. 2013.  
1284 Mice selectively bred for high voluntary wheel running have larger midbrains:  
1285 support for the mosaic model of brain evolution. *J. Exp. Biol.* 216:515–523.

1286 Lerner IM, Dempster ER. 1951. Attenuation of genetic progress under continued  
1287 selection in poultry. *Heredity* 5:75–94.

1288 Lightfoot JT, De Geus EJC, Booth FW, Bray MS, Den Hoed M, Kaprio J, Kelly SA,  
1289 Pomp D, Saul MC, Thomis MA, et al. 2018. Biological/genetic regulation of  
1290 physical activity level: Consensus from GenBioPAC. *Med. Sci. Sports Exerc.*  
1291 50:863–873.

1292 Lillie M, Honaker CF, Siegel PB, Carlborg Ö. 2019. Bidirectional selection for body  
1293 weight on standing genetic variation in a chicken model.  
1294 *Genes|Genomes|Genetics*:g3.400038.2019.

1295 Long A, Liti G, Luptak A, Tenaillon O. 2015. Elucidating the molecular architecture of  
1296 adaptation via evolve and resequence experiments. *Nat. Rev. Genet.* 16:567–  
1297 582.

1298 MacDonald ST, Bamforth SD, Chen C-M, Farthing CR, Franklyn A, Broadbent C,  
1299 Schneider JE, Saga Y, Lewandoski M, Bhattacharya S. 2008. Epiblastic Cited2  
1300 deficiency results in cardiac phenotypic heterogeneity and provides a mechanism  
1301 for haploinsufficiency. *Cardiovasc. Res.* 79:448–457.

1302 MacKay H, Scott CA, Duryea JD, Baker MS, Laritsky E, Elson AE, Garland, Jr. T,  
1303 Fiorotto ML, Chen R, Li Y, et al. 2019. DNA methylation in AgRP neurons  
1304 regulates voluntary exercise behavior in mice. *Nat. Commun.* [Internet] 10.  
1305 Available from: <http://www.nature.com/articles/s41467-019-13339-3>

1306 Malisch JL, Saltzman W, Gomes FR, Rezende EL, Jeske DR, Garland, Jr. T. 2007.  
1307 Baseline and stress-induced plasma corticosterone concentrations of mice  
1308 selectively bred for high voluntary wheel running. *Physiol. Biochem. Zool.*  
1309 80:146–156.

1310 Mathes WF, Nehrenberg DL, Gordon R, Hua K, Garland, Jr. T, Pomp D. 2010.  
1311 Dopaminergic dysregulation in mice selectively bred for excessive exercise or  
1312 obesity. *Behav. Brain Res.* 210:155–163.

1313 Meek TH, Lonquich BP, Hannon RM, Garland, Jr. T. 2009. Endurance capacity of mice  
1314 selectively bred for high voluntary wheel running. *J. Exp. Biol.* 212:2908–2917.

1315 Middleton KM, Goldstein BD, Guduru PR, Waters JF, Kelly SA, Swartz SM, Garland, Jr.  
1316 T. 2010. Variation in within-bone stiffness measured by nanoindentation in mice  
1317 bred for high levels of voluntary wheel running. *J. Anat.* 216:121–131.

1318 Middleton KM, Shubin CE, Moore DC, Carter PA, Garland, Jr. T, Swartz SM. 2008. The  
1319 relative importance of genetics and phenotypic plasticity in dictating bone  
1320 morphology and mechanics in aged mice: Evidence from an artificial selection  
1321 experiment. *Zoology* 111:135–147.

1322 Nguyen QAT, Hillis D, Katada S, Harris T, Pontrello C, Garland, Jr. T, Haga-Yamanaka  
1323 S. 2020. Coadaptation of the chemosensory system with voluntary exercise  
1324 behavior in mice. *PLOS ONE* 15:e0241758.

1325 Porto A, Schmelter R, VandeBerg JL, Marroig G, Cheverud JM. 2016. Evolution of the  
1326 genotype-to-phenotype map and the cost of pleiotropy in mammals. *Genetics*  
1327 204:1601–1612.

1328 Reeve JP. 2000. Predicting long-term response to selection. *Genet. Res.* 75:83–94.

1329 Reeve JP, Fairbairn DJ. 2001. Predicting the evolution of sexual size dimorphism:  
1330 Predicting the evolution of SSD. *J. Evol. Biol.* 14:244–254.

1331 Rhodes JS, Gammie SC, Garland, Jr. T. 2005. Neurobiology of mice selected for high  
1332 voluntary wheel-running activity. *Integr. Comp. Biol.* 45:438–455.

1333 Rhodes JS, Hosack GR, Girard I, Kelley AE, Mitchell GS, Garland, Jr. T. 2001.  
1334 Differential sensitivity to acute administration of cocaine, GBR 12909, and  
1335 fluoxetine in mice selectively bred for hyperactive wheel-running behavior.  
1336 *Psychopharmacology (Berl.)* 158:120–131.

1337 Rhodes JS, van Praag H, Jeffrey S, Girard I, Mitchell GS, Garland, Jr. T, Gage FH.  
1338 2003. Exercise increases hippocampal neurogenesis to high levels but does not  
1339 improve spatial learning in mice bred for increased voluntary wheel running.  
1340 *Behav. Neurosci.* 117:1006–1016.

1341 Roberts RC. 1966. The limits to artificial selection for body weight in the mouse II. The  
1342 genetic nature of the limits. *Genet. Res.* 8:361–375.

1343 Rose MR, Passananti HB, Chippindale AK, Phelan JP, Matos M, Teotónio H, Mueller  
1344 LD. 2005. The effects of evolution are local: evidence from experimental  
1345 evolution in drosophila. *Integr. Comp. Biol.* 45:486–491.

1346 Schlötterer C, Kofler R, Versace E, Tobler R, Franssen SU. 2015. Combining  
1347 experimental evolution with next-generation sequencing: a powerful tool to study  
1348 adaptation from standing genetic variation. *Heredity* 114:431–440.

1349 Schmidt S, Gawlik V, Höltner SM, Augustin R, Scheepers A, Behrens M, Wurst W,  
1350 Gailus-Durner V, Fuchs H, de Angelis MH, et al. 2008. Deletion of glucose  
1351 transporter GLUT8 in mice increases locomotor activity. *Behav. Genet.* 38:396–  
1352 406.

1353 Schwartz NL, Patel BA, Garland, Jr. T, Horner AM. 2018. Effects of selective breeding  
1354 for high voluntary wheel-running behavior on femoral nutrient canal size and  
1355 abundance in house mice. *J. Anat.* 233:193–203.

1356 Sella G, Barton NH. 2019. Thinking about the evolution of complex traits in the era of  
1357 genome-wide association studies. *Annu. Rev. Genomics Hum. Genet.* 20:461–  
1358 493.

1359 Stephan W. 2016. Signatures of positive selection: from selective sweeps at individual  
1360 loci to subtle allele frequency changes in polygenic adaptation. *Mol. Ecol.* 25:79–  
1361 88.

1362 Stratford T, Logan C, Zile M, Maden M. 1999. Abnormal anteroposterior and  
1363 dorsoventral patterning of the limb bud in the absence of retinoids. *Mech. Dev.*  
1364 81:115–125.

1365 Swallow JG, Carter PA, Garland, Jr. T. 1998. Artificial selection for increased wheel-  
1366 running behavior in house mice. *Behav. Genet.* 28:227–237.

1367 Swallow JG, Garland, Jr. T, Carter PA, Zhan W-Z, Sieck GC. 1998. Effects of voluntary  
1368 activity and genetic selection on aerobic capacity in house mice (*Mus*  
1369 *domesticus*). *J. Appl. Physiol.* 84:69–76.

1370 Swallow JG, Hayes JP, Koteja P, Garland, Jr. T. 2009. Selection experiments and  
1371 experimental evolution of performance and physiology. In: Experimental  
1372 evolution: concepts, methods, and applications of selection experiments.  
1373 Berkeley: University of California Press. p. 301–351.

1374 Syme DA, Evashuk K, Grintuch B, Rezende EL, Garland, Jr. T. 2005. Contractile  
1375 abilities of normal and “mini” triceps surae muscles from mice ( *Mus domesticus* )  
1376 selectively bred for high voluntary wheel running. *J. Appl. Physiol.* 99:1308–1316.

1377 Thompson Z, Argueta D, Garland, Jr. T, DiPatrizio N. 2017. Circulating levels of  
1378 endocannabinoids respond acutely to voluntary exercise, are altered in mice  
1379 selectively bred for high voluntary wheel running, and differ between the sexes.  
1380 *Physiol. Behav.* 170:141–150.

1381 Travisano M, Shaw RG. 2013. Lost in the map. *Evolution* 67:305–314.

1382 Val P, Martinez-Barbera J-P, Swain A. 2007. Adrenal development is initiated by Cited2  
1383 and Wt1 through modulation of Sf-1 dosage. *Development* 134:2349–2358.

1384 Verdelis K, Ling Y, Sreenath T, Haruyama N, MacDougall M, Van Der Meulen MCH,  
1385 Lukashova L, Spevak L, Kulkarni AB, Boskey AL. 2008. DSPP effects on in vivo  
1386 bone mineralization. *Bone* 43:983–990.

1387 Wallace IJ, Garland, Jr. T. 2016. Mobility as an emergent property of biological  
1388 organization: Insights from experimental evolution: Mobility and biological  
1389 organization. *Evol. Anthropol. Issues News Rev.* 25:98–104.

1390 Wallace IJ, Middleton KM, Lublinsky S, Kelly SA, Judex S, Garland, Jr. T, Demes B.  
1391 2010. Functional significance of genetic variation underlying limb bone  
1392 diaphyseal structure. *Am. J. Phys. Anthropol.* 143:21–30.

1393 Wallace IJ, Tommasini SM, Judex S, Garland T, Demes B. 2012. Genetic variations and  
1394 physical activity as determinants of limb bone morphology: An experimental  
1395 approach using a mouse model. *Am. J. Phys. Anthropol.* 148:24–35.

1396 Wang Z, Emmerich A, Pillon NJ, Moore T, Hemerich D, Cornelis MC, Mazzaferro E,  
1397 Broos S, Ahluwalia TS, Bartz TM, et al. 2022. Genome-wide association  
1398 analyses of physical activity and sedentary behavior provide insights into  
1399 underlying mechanisms and roles in disease prevention. *Nat. Genet.* 54:1332–  
1400 1344.

1401 Waters RP, Pringle RB, Forster GL, Renner KJ, Malisch JL, Garland, Jr. T, Swallow JG.  
1402 2013. Selection for increased voluntary wheel-running affects behavior and brain  
1403 monoamines in mice. *Brain Res.* 1508:9–22.

1404 Wood AR, The Electronic Medical Records and Genomics (eMERGE) Consortium, The  
1405 MIGen Consortium, The PAGE Consortium, The LifeLines Cohort Study, Esko T,  
1406 Yang J, Vedantam S, Pers TH, Gustafsson S, et al. 2014. Defining the role of  
1407 common variation in the genomic and biological architecture of adult human  
1408 height. *Nat. Genet.* 46:1173–1186.

1409 Xie Y, Pan W, Khodursky AB. 2005. A note on using permutation-based false discovery  
1410 rate estimates to compare different analysis methods for microarray data.  
1411 *Bioinformatics* 21:4280–4288.

1412 Xu S, Garland T. 2017. A mixed model approach to genome-wide association studies  
1413 for selection signatures, with application to mice bred for voluntary exercise  
1414 behavior. *Genetics* 207:785–799.

1415 Zamer WE, Scheiner SM. 2014. A conceptual framework for organismal biology: linking  
1416 theories, models, and data. *Integr. Comp. Biol.* 54:736–756.

1417 Zhang M, Han Y, Liu J, Liu L, Zheng L, Chen Y, Xia R, Yao D, Cai X, Xu X. 2020.  
1418 Rbm24 modulates adult skeletal muscle regeneration via regulation of alternative  
1419 splicing. *Theranostics* 10:11159–11177.

1420 Zinski AL, Carrion S, Michal JJ, Gartstein MA, Quock RM, Davis JF, Jiang Z. 2021.  
1421 Genome-to-phenome research in rats: progress and perspectives. *Int. J. Biol.  
1422 Sci.* 17:119–133.

1423

1424

1425 **Supporting Information**

1426 **S1 Fig. Simulation power results by effect size and generation.** Power (Y-axis) of  
1427 different effect sizes at different generations (X-axis). Effect Size - Color: 204.8 - brown,  
1428 102.4 - red, 51.2 - orange, 25.6 - yellow, 12.8 - dark green, 6.4 - light green, 3.2 - dark  
1429 blue, 1.6 - light blue, 0.8 - dark purple, 0.4 - light purple

1430

1431 **S1 File. Regularized and windowed regularized F-test (WRT).** Description of  
1432 methodology and rationale.

1433

1434 **S2 File. Parameters and seeds for constraint simulations.** List of the parameters  
1435 and seeds used for the simulations with and without constraints.

1436

1437 **S1 Table. Effect size and Type I error rates.** Includes error rates and means for  
1438 different effect sizes and sample sizes.

1439

1440 **S2 Table. Power to detect differentiation between HR and C lines in relation to  
1441 effect size of locus and generation under two simulation models.** Includes effect  
1442 sizes, sample sizes, revolutions when homozygous, and power and mean for each  
1443 generation and model.

1444

1445 **S3 Table. Differentiated regions identified at generation 22 (FDR = 0.01).** Includes  
1446 chromosomal location, size of region, most statistically significant base pair p-value, and  
1447 position of this base pair.

1448

1449 **S4 Table. Genes included in “strict” culling regions at generation 22.**

1450

1451 **S5 Table. Gene ontology results for generation 22 “strict” culling genes.** Includes  
1452 GO terms, fold enrichment, and raw p-values.

1453