1    **Title: Decreased adaptation at human disease genes as a possible consequence of**

2    **interference between advantageous and deleterious variants**

3

4    **Authors**

5    Chenlu Di[1], Diego Salazar Tortosa[1], M. Elise Lauterbur[1] and David Enard[1]

6

7    **Affiliation**

8    [1] University of Arizona Department of Ecology and Evolutionary Biology, Tucson, Arizona,
9    USA.
10

11    **Corresponding author**

12    David Enard, denard@email.arizona.edu

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

**Abstract**

**Advances in genome sequencing have dramatically improved our understanding of the genetic basis of human diseases, and thousands of human genes have been associated with different diseases. Despite our expanding knowledge of gene-disease associations, and despite the medical importance of disease genes, their evolution has not been thoroughly studied across diverse human populations. In particular, recent genomic adaptation at disease genes has not been well characterized, even though multiple evolutionary processes are expected to connect disease and adaptation at the gene level. Understanding the relationship between disease and adaptation at the gene level in the human genome is severely hampered by the fact that we don't even know whether disease genes have experienced more, less, or as much adaptation as non-disease genes during recent human evolution. Here, we compare the rate of strong recent adaptation in the form of selective sweeps between disease genes and non-disease genes across 26 distinct human populations from the 1,000 Genomes Project. We find that disease genes have experienced far less selective sweeps compared to non-disease genes during recent human evolution. This sweep deficit at disease genes is particularly visible in Africa, and less visible in East Asia or Europe, likely due to more intense genetic drift in the latter populations creating more spurious selective sweeps signals. Investigating further the possible causes of the sweep deficit at disease genes, we find that this deficit is very strong at disease genes with both low recombination rates and with high numbers of associated disease variants, but is inexistant at disease genes with higher recombination rates or lower numbers of associated disease variants. Because recessive deleterious variants have the ability to interfere with adaptive ones, these observations strongly suggest that adaptation has been slowed down by the presence of interfering recessive deleterious variants at disease genes. These results clarify the evolutionary relationship between disease genes and recent genomic adaptation, and suggest that disease genes suffer not only from a higher load of segregating deleterious mutations, but also an inability to adapt as much, and/or as fast as the rest of the genome.**

63    **Keywords:** adaptation, human disease, Hill–Robertson interference, recessive deleterious

64    variants, selective sweeps, environmental changes

65    **Introduction**

66    Advances in genome sequencing have dramatically improved our understanding of the genetic

67    basis of human diseases, and thousands of human genes have been associated with different

68    diseases (Amberger et al., 2019; Piñero et al., 2020). Despite our expanding knowledge of gene-

69    disease associations, and despite the fact that multiple evolutionary processes might connect

70    disease and genomic adaptation at the gene level, these connections are yet to be studied.

71    Different evolutionary processes have the potential to make the occurrence of disease genes and

72    adaptation not independent from each other in the human genome. For instance, hitchhiking of

73    deleterious mutations linked to advantageous mutations might increase the risk of disease-

74    causing variants at genes subjected to past directional adaptation. Disease genes might then

75    appear to have experienced more adaptation than non-disease genes if this specific process was

76    sufficiently widespread. Conversely, higher evolutionary constraint, and higher pleiotropy might

77    reduce adaptation at disease genes compared to genes not involved in diseases (Otto, 2004).

78    There is currently considerable uncertainty about how any of these non-exclusive evolutionary

79    processes, or other processes, might have influenced adaptation at disease genes. It is even not

80    well-known whether human non-infectious disease genes have similar, higher or lower levels of

81    adaptation in human populations compared to genes not involved in diseases. Comparing levels

82    of adaptation between disease genes and non-disease genes is a first important step toward better

83    understanding the evolutionary relationship between non-infectious diseases and genomic

84    adaptation.

85

86        Multiple recent studies comparing evolutionary patterns between human disease and non-

87    disease genes have found that disease genes are more constrained and evolve more slowly (lower

88    ratio of nonsynonymous to synonymous substitution rate, dN/dS, in disease genes) (Blekhman et

89    al., 2008; Park et al., 2012; Spataro et al., 2017), An older comparison by Smith and Eyre-Waler

90    (2003) found that disease genes evolve faster than non-disease genes (higher dN/dS), but we note

91    that the sample of disease genes used at the time was very limited.
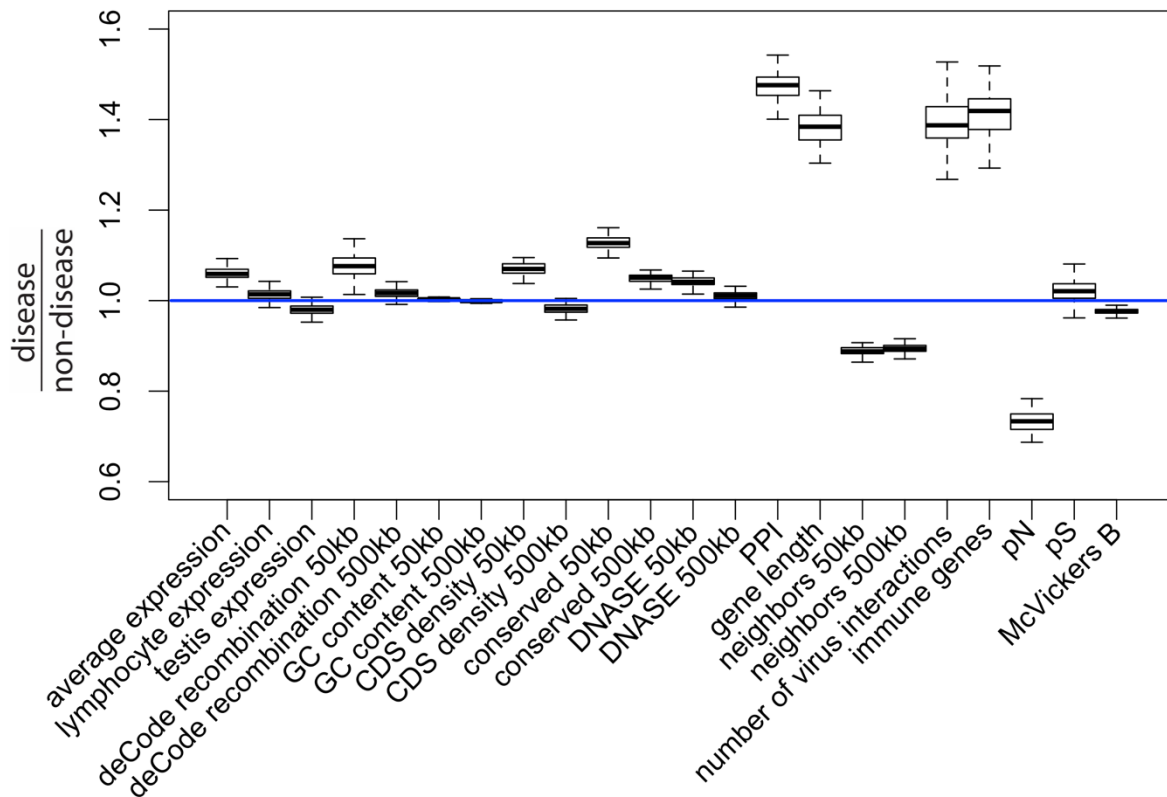
92

93   The significant increase of the number of known disease genes since these studies were

94 completed makes it important to update the comparison of evolutionary patterns at disease and

95 non-disease genes. More critically however, past studies all have in common an important

96 limitation that justifies comparing disease genes and non-disease genes again.  Disease and non-

97 disease genes may differ by more than just the fact that they have been associated with disease or

98 not. Disease and non-disease genes may also differ in many other factors other than their disease

99 status. Such factors can be a problem when comparing adaptation in disease genes and non-

100 disease genes, because they, instead of the disease status itself, could explain differences in

101 adaptation. For example, disease genes tend to be more highly expressed than non-disease genes

102 (Spataro et al., 2017) (Figure 1). If higher expression happens to be associated with more

103 adaptation in general, one might detect more adaptation in disease genes in a way that has

104 nothing to do with disease, and just reflects their higher levels of expression. Many other factors

105 may also be important. For example, immune genes, which often adapt in response to infectious

106 pathogens, may further complicate comparisons if they are represented in unequal proportions

107 between non-infectious disease and non-disease genes. Comparing genomic adaptation in disease

108 and non-disease genes thus requires careful consideration of confounding factors.

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125



126

**Figure 1. Potential confounding factors in disease versus non-disease genes.**
Each potential confounding factor is detailed in the Methods. For each confounding factor, the boxplot shows on the y-axis the ratio of the average factor value for disease genes, divided by the average factor value for non-disease genes. The boxplot error bars are obtained by calculating the ratio 1,000 times, each time by randomly sampling as many non-disease genes as there are disease genes.

133

134

Among other confounding factors, it is particularly important to take into account evolutionary constraint, i.e the level of purifying selection experienced by different genes. A common intuition is that disease genes may exhibit less adaptation because they are more constrained (Blekhman et al., 2008), leaving less mutational space for adaptation to happen in the first place. Less adaptation at disease genes might thus represent a trivial consequence of varying constraint between genes (Kim et al., 2007), which says little about a specific connection

5

141  between disease and adaptation. In the same vein, one might expect disease genes to be

142  associated with higher mutation rates, and more frequent adaptation to follow as a trivial

143  consequence of elevated mutation rates. Whether disease genes experience higher mutation rates

144  is however still an open question (Osada et al., 2009; Eyre-Walker and Eyre-Walker, 2014). In

145  any case, focusing specifically on disease and adaptation requires controlling for confounders

146  such as constraint and mutation rate (see Methods, Results and Figure 1 for a complete list of

147  confounders accounted for in this analysis).

148

149      A specific evolutionary relationship may exist between adaptation and disease beyond the

150  simple effect of constraint, mutation rate or other confounders. In an evolutionary context, once

151  constraint and other confounding factors have been accounted for, we can imagine three potential

152  scenarios for the comparison of adaptation between disease and non-disease genes. Under

153  scenario 1, any potential difference in adaptation between disease and non-disease genes is

154  entirely due to differences in constraint and other confounding factors. Under this scenario, there

155  is no further evolutionary process linking disease and adaptation together. Therefore, there is no

156  difference in adaptation between disease and non-disease genes once confounding factors have

157  been accounted for.

158

159      Under scenario 2, disease genes have more adaptation than non-disease genes. For

160  example, as already mentioned above, deleterious mutations can hitchhike together with adaptive

161  mutations to high frequencies in human populations (Birky and Walsh, 1988; Barreiro and

162  Quintana-Murci, 2010; Chun and Fay, 2011). Other, less well established, cases can be imagined

163  where past adaptation decreased the robustness of a specific gene, and subsequent mutations

164  become more likely to be associated with diseases (Xu and Zhang, 2014). Scenario 2 thus favors

165  a relationship between adaptation and disease, where past adaptation precedes and influences the

166  likelihood of a gene being associated with disease.

167      Under scenario 3, disease genes have less adaptation than non-disease genes even after

168  accounting for confounding factors such as evolutionary constraint. Such a scenario might occur

169  for example if disease genes happen to be genes that can be sensitive to changes in the

170  environment, with a fitness optimum that can change over time, but where adaptation has not

171  occurred yet to catch up with the new optimum. Such an adaptation lag (or lag load, to reuse the

6

172     terminology introduced by J. Maynard-Smith (1976)) may occur for example if higher pleiotropy

173     at disease genes (Ittisoponpisan et al., 2017) makes it less likely for new mutations to be

174     advantageous (Otto, 2004) (in addition to increasing the level of constraint already accounted for

175     as a confounding factor).  Such an adaptation lag, with genes further away from their optimum,

176     might make such genes more prone to accumulate disease variants that fall too far from the

177     "normal" functioning range around the optimum. An adaptation lag may also occur if deleterious

178     mutations interfere with and slow down adaptation at disease genes more than at non-disease

179     genes (Assaf et al., 2015; Hill and Robertson, 1966).

180         Even though uncovering the underlying evolutionary processes that govern the

181     relationship between disease and adaptation will take a lot more work than the present analysis, it

182     is important to find first which scenario is the most likely to be true, i.e whether disease genes

183     have as much, more, or less adaptation than non-disease genes. Finding out which out of the

184     three possible scenarios is true may give a preliminary basis to further hypothesize which

185     evolutionary processes are more likely to dominate the relationship between disease and

186     adaptation genome-wide.

187

188         Here, we compare recent adaptation in mendelian disease and non-disease genes in order

189     to disentangle the connections between adaptation and disease. We specifically compare the

190     abundance of recent selective sweeps signals, where hitchhiking has raised haplotypes that carry

191     an advantageous variant to higher frequencies (Smith and Haigh, 1974). Note that this means that

192     we can only compare adaptation at specific loci between disease and non-disease genes that was

193     strong enough to induce hitchhiking, hence we do not take into account polygenic adaptation

194     distributed across a large number of loci that did not leave any hitchhiking signals (see

195     Discussion). As mentioned above, confounding factors may affect the comparison between

196     disease and non-disease genes. In contrast with previous studies, we systematically control for a

197     large number of confounding factors when comparing recent adaptation in human disease and

198     non-disease genes, including evolutionary constraint, mutation rate, recombination rate, the

199     proportion of immune or virus-interacting genes, etc. (please refer to Methods for a full list of the

200     confounding factors included). In addition to controlling for a large number of confounding

201     factors, we estimate false positive risks (FPR) for our comparison pipeline that fully take into

202     account the implications of controlling for many factors (see Methods and Results).

203    As a list of disease genes to test, we curate human mendelian non-infectious disease

204    genes based on annotations in the DisgeNet and OMIM databases (Methods). We focus on

205    mendelian disease genes rather than all disease genes including complex disease associations,

206    because different evolutionary patterns can be expected between mendelian and complex disease

207    genes based on previous studies (Blekhman et al., 2008; Quintana-murci, 2016; Spataro et al.,

208    2017). In total, we compare 4,215 mendelian disease genes with non-disease genes in the human

209    genome. In agreement with scenario 3, we find a strong deficit of selective sweeps at disease

210    genes compared to non-disease genes. We further test multiple potential explanations for this

211    deficit, and find that higher pleiotropy at disease genes is unlikely to explain the less frequent

212    occurrence of sweeps. In contrast, we find that the sweep deficit at disease genes strongly

213    depends on recombination and the number of known disease variants at given disease genes.

214    This suggests that segregating deleterious mutations at disease genes might interfere with, and

215    slow down genetically linked adaptive variants enough to produce the observed lack of sweeps at

216    disease genes.

217

218    **Results**

219

220    **Controlling for confounding factors with a bootstrap test**

221    To compare disease and non-disease genes, we first ask which potential confounding factors

222    differ between the two groups of genes. As expected, multiple measures of selective constraint

223    are significantly higher in disease compared to non-disease genes. As a measure of long-term

224    constraint, the density of conserved elements across mammals is slightly higher at disease genes

225    compared to non-disease genes (Figure 1: conserved 50kb, conserved 500kb; Methods).

226    As a measure of more recent constraint, we contrast pS, the average proportion of variable

227    synonymous sites, with pN, the average proportion of variable nonsynonymous sites (Figure 1;

228    Methods). If the coding sequences of disease genes are more constrained, we expect a drop of pN

229    at disease genes, but no such drop of pS at neutral synonymous sites. Accordingly, pN is lower at

230    disease compared to non-disease genes, while pS is very similar between the two categories of

231    genes (Figure 1). Therefore, selective constraint was stronger in the coding sequences of disease

232    genes during recent human evolution.

233    As another measure of recent constraint, we also use McVicker's B estimator of background

234    selection (McVicker et al., 2009). The amount of background selection at a locus can be used as

235    a proxy for recent constraint, since it depends on the number of deleterious mutations that were

236    recently removed at this locus. The lower B, the more background selection there is at a specific

237    locus. In line with higher recent constraint at disease genes, B is slightly, but significantly lower

238    at disease genes (Figure 1; Methods). Overall, we find evidence of higher constraint at disease

239    genes.

240

241    In addition to constraint, mutation rate could represent an important confounder. The proportion

242    of variable neutral synonymous sites pS can be used to compare mutation rates, since the number

243    of variable sites is proportional to the mutation rate under neutrality. As mentioned already, pS is

244    very similar at disease and non-disease genes (Figure 1), suggesting that mutation rates are

245    similar at disease and non-disease genes. This is further supported by the fact that multiple

246    factors that could affect the mutation rate such as GC content or recombination are also similar at

247    disease and non-disease genes (Figure 1; Methods). Aside from mutation rate and constraint,

248    multiple other factors that could affect adaptation differ between disease and non-disease genes,

249    notably including the proportion of genes that interact with viruses, the proportion of immune

250    genes, or the number of protein-protein interactions (PPIs) in the human PPIs network. All these

251    factors have been shown to affect adaptation (Methods), further showing the necessity to control

252    for confounding factors when comparing adaptation at disease and non-disease genes.

253

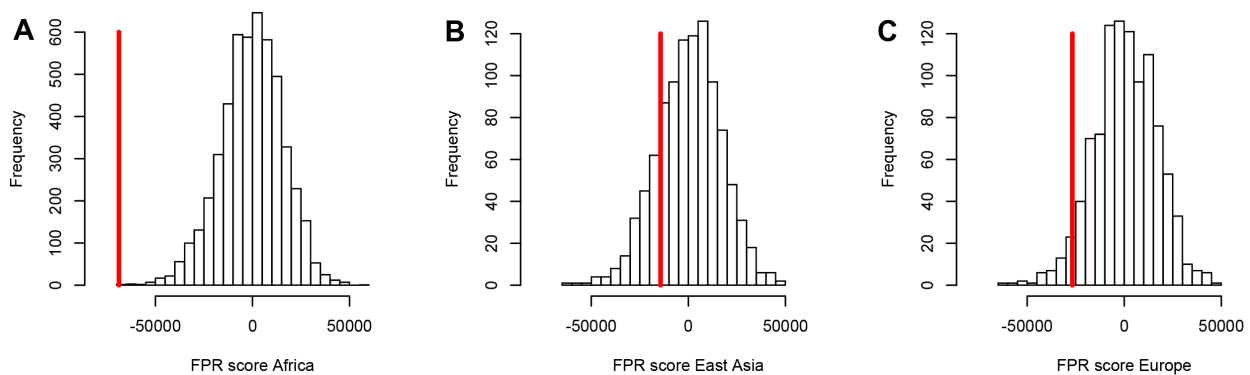254    **Less sweeps at disease genes**

255        For our comparison of disease and non-disease genes, we measure recent adaptation

256    around human protein coding genes (Methods) using the integrated haplotype score (iHS,

257    (Voight et al., 2006)) and the number of Segregating sites by Length ($nS_L$, (Ferrer-Admetlla et

258    al., 2014)) in 26 populations (The 1000 Genomes Project Consortium, 2015) (Methods). The iHS

259    and $nS_L$ statistics are both sensitive to recent incomplete sweeps, and have the advantage over

260    other sweep statistics of being insensitive to the confounding effect of background selection

261    (Enard et al., 2014; Schrider, 2020). To evaluate the prevalence of sweeps at disease genes

262    relative to non-disease genes, we do not use the classic outlier approach, and instead used a

263    previously described, more versatile approach based on block-randomized genomes to estimate

264    unbiased false positive risks for whole enrichment curves (Figure 2) (Enard and Petrov, 2020).

265    We first rank genes based on the average iHS or $nS_L$ in genomic windows centered on genes

266    (Methods), from the top-ranking genes with the strongest sweep signals to the genes with the

267    weakest signals. We then slide a rank threshold from a high rank value to a low rank value (from

268    top 5,000 to top 10, x-axis on Figure 2). For each rank threshold, we estimate the sweep

269    enrichment (or deficit) at disease relative to non-disease genes (Figure 2, y-axis). For example,

270    for rank threshold 200, the relative enrichment (or deficit) is the number of disease genes in the

271    top 200 ranking genes, divided by the number of control non-disease genes in the top 200. By

272    sliding the rank threshold, we estimate a whole enrichment curve that is not only sensitive to the

273    strongest sweeps but also to weaker sweeps signals (for example using the top 5,000 threshold;

274    Figure 2). Using block-randomized genomes (Methods), we can then estimate an unbiased false

275    positive risk (FPR) for the whole enrichment curve. This strategy makes less assumptions on the

276    expected strength of selective sweeps. The approach also makes it possible to estimate a single

277    false positive risk based on the cumulated enrichment (or deficit) over multiple whole

278    enrichment curves (Methods). Here, we estimate a single false positive risk for both iHS and $nS_L$

279    curves considered together, and also for multiple window sizes to measure average iHS and $nS_L$

280    (from 50kb to 1Mb, Methods).

281

282    To control for confounding factors (Figure 1), we compare sweep signals at disease genes with

283    control non-disease genes that were chosen by a bootstrap test (Castellano et al., 2019; Enard and

284    Petrov, 2020) because they match disease genes in terms of confounding factor values

285    (Methods). Furthermore, control non-disease genes are chosen far from disease genes (>300kb;

286    Methods). We do this to avoid choosing as controls non-disease genes that are too close to

287    disease genes and thus likely to have the same sweep profile (especially in the case of large

288    sweeps potentially overlapping both neighboring disease and non-disease genes). This, together

289    with the large number of confounding factors that we match, tends to limit the pool of possible

290    control genes (Methods). The statistical impact of a limited control pool is however fully taken

291    into account by the estimation of a FPR with block-randomized genomes (Methods).

292

293    Because they have experienced different demographic histories, we test different human

294    populations from distinct continents separately. Specifically, we test African populations, East

10

295    Asian populations and European populations from the 1,000 Genomes Project phase 3 (The 1000

296    Genomes Project Consortium, 2015). At this stage we must consider the fact that most gene-

297    disease associations in our dataset were likely discovered in European cohorts. Because disease

298    genes in Europe may not always be disease genes in other populations, we cannot exclude the

299    possibility that a sweep enrichment or a sweep deficit might be more pronounced in Europe,

300    unless the evolutionary processes that make a gene more likely to be a disease gene predated the

301    split of different human populations. Conversely, one might expect distinct selective patterns

302    between disease and non-disease genes to be more visible in Africa. Indeed, more intense drift,

303    due to the more severe bottlenecks experienced by ancestral Eurasian populations (The 1000

304    Genomes Project Consortium, 2015), is expected to dilute true selective patterns among false

305    positive signals more in Europe and East Asia, by creating a higher base level of drift noise.
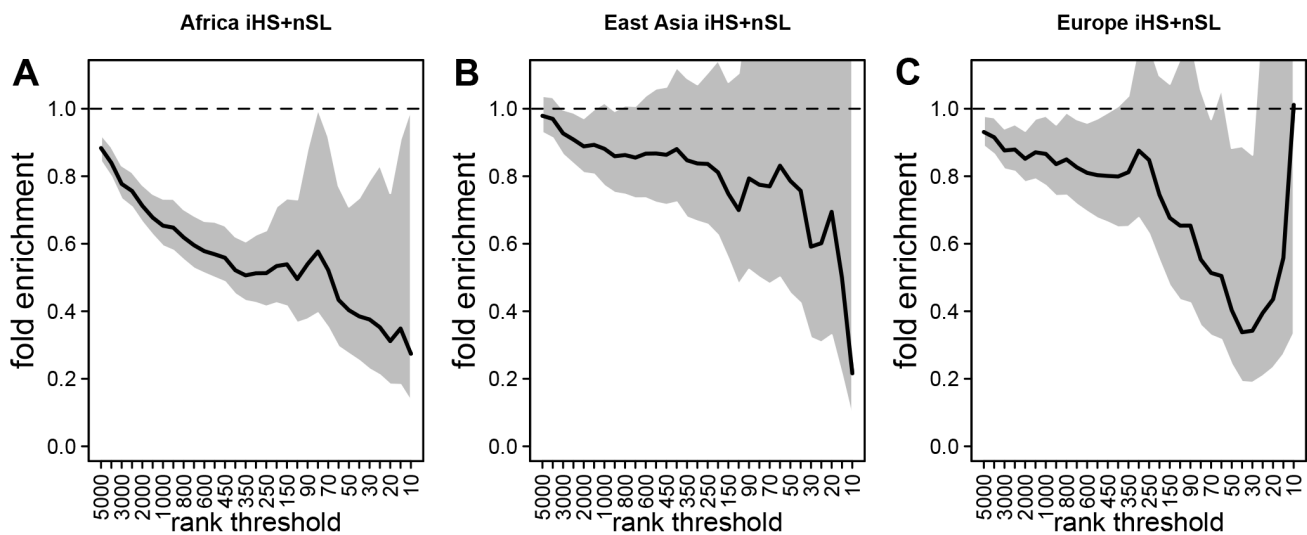
306



308    **Figure 2. A stronger sweep deficit at disease genes in Africa than in East Asia and Europe.**
309    The figure shows the observed sweep enrichment/deficit score used to measure the false positive
310    risk (FPR) in the real genome (red line), compared to the expected null distribution of the score
311    estimated with block-randomized genomes (5,000 block-randomized genomes in Africa, 1,000 in
312    East Asia and Europe; Methods). The FPR score is based on summing the difference between the
313    number of genes in sweeps at disease genes and the number of genes in sweeps in control genes,
314    over both iHS and $nS_L$, and different window sizes (Methods). A) FPR score in Africa, estimated
315    summing over the ESN, GWD, LWK, MSL and YRI populations from the 1,000 Genomes
316    Project. B) FPR score in East Asia, estimated summing over the CDX, CHB, CHS, JPT and
317    KHV populations. C) FPR score in Europe, summing over the CEU, FIN, GBR, IBS and TSI
318    populations.

319

320

321    Using both iHS and $nS_L$ sweep signals, we find a strong depletion in sweep signals at disease

322    genes, especially in Africa with a low false positive risk (FPR=3.10-4 vs. 0.18 in East Asia and

323    0.05 in Europe, Figure 2A, B and C respectively; Methods). Note that this FPR takes the

11

324     clustering of multiple genes in the same sweeps into account (Enard and Petrov, 2020). A

325     stronger depletion in Africa suggests that the evolutionary processes linking disease and

326     adaptation at the gene level predate the split of African and European populations, given that

327     most gene-disease associations studies involved European cohorts. The stronger depletion in

328     Africa also suggests that the same pattern might be present outside of Africa, but more hidden by

329     genetic drift noise. It might indeed be harder to distinguish a deficit of true sweep signals at

330     disease genes if it is swamped by an elevated level of false sweep signals occurring at random in

331     the genome, due to more intense drift. Figure 3A, B and C show the sweep deficit curves at

332     disease genes compared to control non-disease genes in Africa, East Asia and Europe,

333     respectively.

334

335



336

337     **Figure 3. Deficit of iHS and $nS_L$ sweep signals at disease genes.**
338     The figure shows the averaged whole enrichment curves and their averaged confidence intervals
339     from the bootstrap test, averaged over both iHS and $nS_L$ sweep ranks, and over all the
340     populations from each continent (Methods). The y-axis represents the relative sweep enrichment
341     at disease genes, calculated as the number of disease genes in putative sweeps, divided by the
342     number of control non-disease genes in putative sweeps. The gray areas are the 95% confidence
343     interval for this ratio. The number of genes in putative sweeps is measured for varying sweep
344     rank thresholds. For example, at the top 100 rank threshold, the relative enrichment is the
345     number of disease genes within the top 100 genes with the strongest sweep signals (either
346     according to iHS or $nS_L$), divided by the number of control non-disease genes within the top 100
347     genes with the strongest sweep signals. We use genes ranked by iHS or $nS_L$ using 200kb
348     windows, since 200kb is the intermediate size of all the window sizes we use (50kb, for the
349     smallest, 1000kb for the largest; see Methods). A) Africa, average over the ESN, GWD, LWK,

12

350  MSL and YRI populations from the 1,000 Genomes Project. B) East Asia, average over the
351  CDX, CHB, CHS, JPT and KHV populations. C) Europe, average over the CEU, FIN, GBR, IBS
352  and TSI populations.
353

354  Notably, the stronger depletion observed in Africa likely excludes the possibility that it could be

355  mostly due to a technical artifact, where sweeps themselves might make it harder to identify

356  disease genes in the first place. Sweeps increase linkage disequilibrium (LD) in a way that could

357  make it more difficult to assign a disease to a single gene in regions of the genome with high LD

358  and multiple genes genetically linked to a disease variant. This could result in a depletion of

359  sweeps at monogenic disease genes, simply because disease genes are less well annotated in

360  regions of high LD. However, if this was the case, because most disease gene were identified in

361  Europe, we would expect such an artifact to deplete sweeps at disease genes primarily in Europe,

362  not in Africa. This artifact is also very unlikely due to the fact that recombination rates are

363  similar between disease and non-disease genes (Figure 1). Overall, these results support the third

364  scenario where evolutionary processes decrease adaptation at disease genes. That said, it is

365  important to note that we only detect a deficit of adaptation strong enough to leave hitchhiking

366  signals. Our results do not imply that the same is true for adaptation that is too polygenic to leave

367  signals detectable with iHS or $nS_L$. Note that the sweep deficit at disease genes in Africa is

368  robust to differences in gene functions between disease and non-disease genes according to a

369  Gene Ontology analysis (Methods) (Gene Ontology Consortium, 2021).

370

371  **A limited role of pleiotropy**

372  A deficit of strong adaptation (strong enough to affect iHS or $nS_L$) raises the question of what

373  creates this deficit at disease genes. Because disease genes tend to be pleiotropic and many

374  disease genes are involved in multiple diseases (see below), pleiotropy is a particularly attractive

375  potential explanation for the lack of sweeps at disease genes. Pleiotropy is defined as the ability

376  for a gene to affect multiple phenotypes. The involvement in multiple phenotypes may make it

377  more difficult for mutations to emerge at pleiotropic genes without any adverse antagonistic

378  effects (Otto, 2004). In addition to the higher selective constraint already accounted for,

379  pleiotropy may thus also make it less likely for advantageous mutations to be advantageous and

380  cause a sweep (Otto, 2004), with the advantage provided by changes at specific phenotypes

381  being mitigated by the adverse effects on other phenotypes.

382     We can test the involvement of pleiotropy with our dataset by comparing sweeps at disease

383     genes involved in multiple diseases, with sweeps at disease genes involved in only one disease.

384     If pleiotropy decreases the rate of sweeps at disease genes, we predict that genes involved in

385     multiple diseases should experience less sweeps than genes involved in only one disease.

386     There are 1221 disease genes in our dataset associated with five or more diseases (five+ disease

387     genes), and 1296 disease genes associated with only one disease according to the CUI (Concept

388     Unique Identifiers) classification provided by DisGeNet (Methods). When comparing the five+

389     disease genes with one disease genes far away (>300 kb as when comparing all disease genes

390     with control non-disease genes), we do not find significantly less iHS and $nS_L$ sweep signals at

391     five+ disease genes in Africa (FPR=0.46). This result makes it unlikely that pleiotropy can

392     explain the sweep deficit at disease genes.

393

394

395     **A possible role of interference of deleterious mutations**

396     With pleiotropy likely having a limited role, we further test other possible explanations for the

397     sweep deficit at disease genes. Another possibility is that adaptation may be limited at disease

398     genes due to deleterious mutations interfering with and slowing down advantageous variants.

399     This process has been mostly studied in haploid species (Peck, 1994; Johnson and Barton, 2002;

400     Jain, 2019). In diploid species including humans, recessive deleterious mutations specifically

401     have been shown to have the ability to slow down, or even stop the frequency increase of

402     advantageous mutations that they are linked with (Assaf et al., 2015; Uricchio et al., 2019).

403     Uricchio et al. (2019) in particular found evidence of decreased protein adaptation in the regions

404     of the human genome with strong background selection and low recombination. The majority of

405     disease variants are recessive (Amberger et al., 2019). Thus, if segregating recessive deleterious

406     mutations are more common at disease genes, starting with the known disease variants

407     themselves, then their interference could in theory explain the sweep deficit that we observe.

408     This is true even despite the fact that we matched disease and control non-disease genes for

409     multiple measures of selective constraint. Indeed, we use measures of selective constraint such as

410     the density of conserved elements or the proportion of variable non-synonymous sites pN

411     (Methods), that are indicative of the amount of deleterious mutations that get ultimately

412     removed, but do not provide any detailed information on either the strength of negative selection,

14

413    or on dominance coefficients. Disease genes and control non-disease genes may have very

414    similar densities of conserved elements and similar pN, and still very different distributions of

415    selection and dominance coefficients of deleterious mutations. Unfortunately, disentangling

416    selection from dominance coefficients is notoriously difficult, because different combinations of

417    selection and dominance coefficients can result in the same patterns of genetic variation (Huber

418    et al., 2018). Although directly comparing the actual total numbers of recessive deleterious

419    mutations at disease and non-disease genes is therefore not possible, we can still use indirect

420    comparison strategies. First, if an interference of deleterious mutations is involved, then this

421    interference is expected to be stronger in low recombination regions of the genome, where more

422    deleterious mutations are likely to be genetically linked to an advantageous mutation. Therefore,

423    we predict that the sweep deficit should be more pronounced when comparing disease and non-

424    disease genes only in low recombination regions of the genome, where the linkage between

425    deleterious and advantageous variants is higher. Conversely, the sweep deficit should be less

426    pronounced in high recombination regions of the genome. Second, if the number of known

427    disease variants at a given disease gene correlates well enough with the total number of

428    segregating recessive deleterious mutations at this disease gene, then we should observe a

429    stronger sweep deficit at disease genes with many known disease variants, compared to disease

430    genes with few known disease variants. Based on these two predictions, the sweep deficit should

431    be particularly strong at disease genes with both many disease variants AND lower

432    recombination. As the number of disease variants for each disease gene, we use the number of

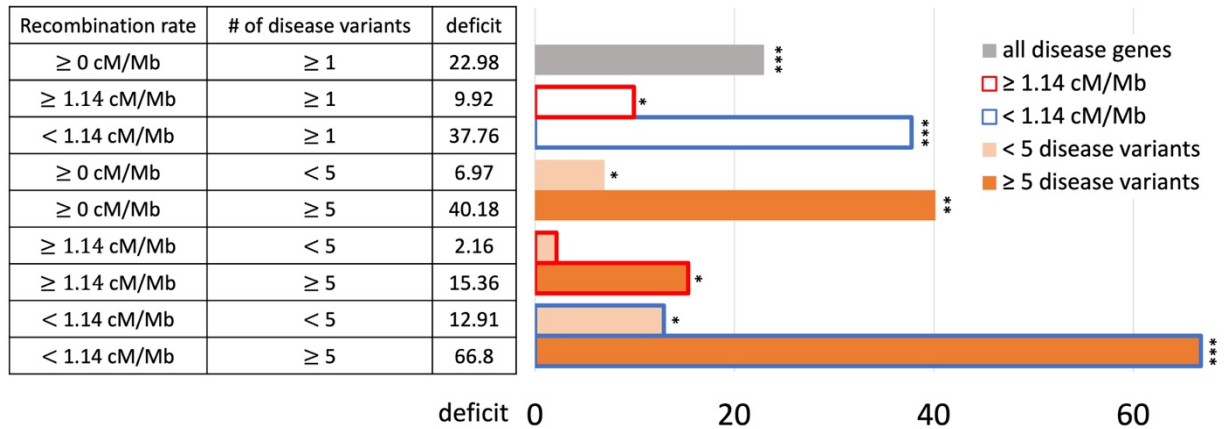433    disease variants as curated by OMIM/UNIPROT (Methods).

434

435    For these comparisons we focus solely on African populations for which we found the strongest

436    sweep deficit (Figure 2). We first compare disease and control non-disease genes both from only

437    regions of the genome with recombination rates lower than the median recombination rate (1.137

438    cM/Mb). In agreement with recombination being involved, we find that the sweep deficit at low

439    recombination disease genes is much more pronounced than the overall sweep deficit found

440    when considering all disease and control non-disease genes regardless of recombination (Figure

441    4, FPR=$2.10^{-4}$). Conversely, the sweep deficit at disease genes compared to non-disease genes is

442    much less pronounced when restricting the comparison to genes with recombination rates higher

443    than the median recombination rate (1.137 cM/Mb), and remains only marginally significant

15

444    (Figure 4, FPR=0.029). This provides evidence that genetic linkage may indeed be involved.

445    Low recombination is however not sufficient on its own to create a sweep deficit, and we further

446    test if the sweep deficit also depends on the number of disease variants at each disease gene. In

447    our dataset, approximately half of all the disease genes have five or more disease variants, and

448    the other half have four or less disease variants (Methods). In further agreement with possible

449    interference of recessive deleterious variants, the sweep deficit is much more pronounced at

450    disease genes with five or more disease variants (Figure 4, FPR=$8.10^{-4}$). The sweep deficit at

451    disease genes with four or less disease variants is barely significant compared to control non-

452    disease genes (Figure 4, FPR=0.032). In addition, disease genes with five or more disease

453    variants, but with recombination higher than the median recombination rate, do not have a strong

454    sweep deficit either (Figure 4, FPR=0.026). A higher number of disease variants alone is thus not

455    enough to explain the sweep deficit. In a similar vein, disease genes with a recombination rate

456    less than the median recombination rate, and with four or less disease variants, do not exhibit a

457    strong sweep deficit (Figure 4, FPR=0.021). This confirms that low recombination alone is not

458    enough to explain the sweep deficit at disease genes. Accordingly, disease genes with both low

459    recombination AND five or more disease variants show the strongest sweep deficit (Figure 4,

460    FPR=$2.10^{-4}$). Disease genes with both high recombination AND less than 5 disease variants show

461    no sweep deficit at all, with a sweep prevalence undistinguishable from control non-disease

462    genes (Figure 4, FPR=0.74). The latter result is important, because it suggests that interference of

463    recessive deleterious variants may be sufficient on its own to explain the whole sweep deficit at

464    disease genes. Both higher linkage and more disease variants seem to be needed to explain the

465    sweep deficit at disease genes. Note that these results are not due to introducing a bias in the

466    overall number of variants by using the number of disease variants, because we always match the

467    level of neutral genetic variation between disease genes and control non-disease genes with pS.

468    The overall level of genetic variation is further matched thanks to pN and thanks to McVicker's

469    B, whose value is directly dependent on the level of genetic variation at a given locus (McVicker

470    et al., 2009).

471

472

473

| Recombination rate | # of disease variants | deficit |
|---|---|---|
| ≥ 0 cM/Mb | ≥ 1 | 22.98 |
| ≥ 1.14 cM/Mb | ≥ 1 | 9.92 |
| < 1.14 cM/Mb | ≥ 1 | 37.76 |
| ≥ 0 cM/Mb | < 5 | 6.97 |
| ≥ 0 cM/Mb | ≥ 5 | 40.18 |
| ≥ 1.14 cM/Mb | < 5 | 2.16 |
| ≥ 1.14 cM/Mb | ≥ 5 | 15.36 |
| < 1.14 cM/Mb | < 5 | 12.91 |
| < 1.14 cM/Mb | ≥ 5 | 66.8 |

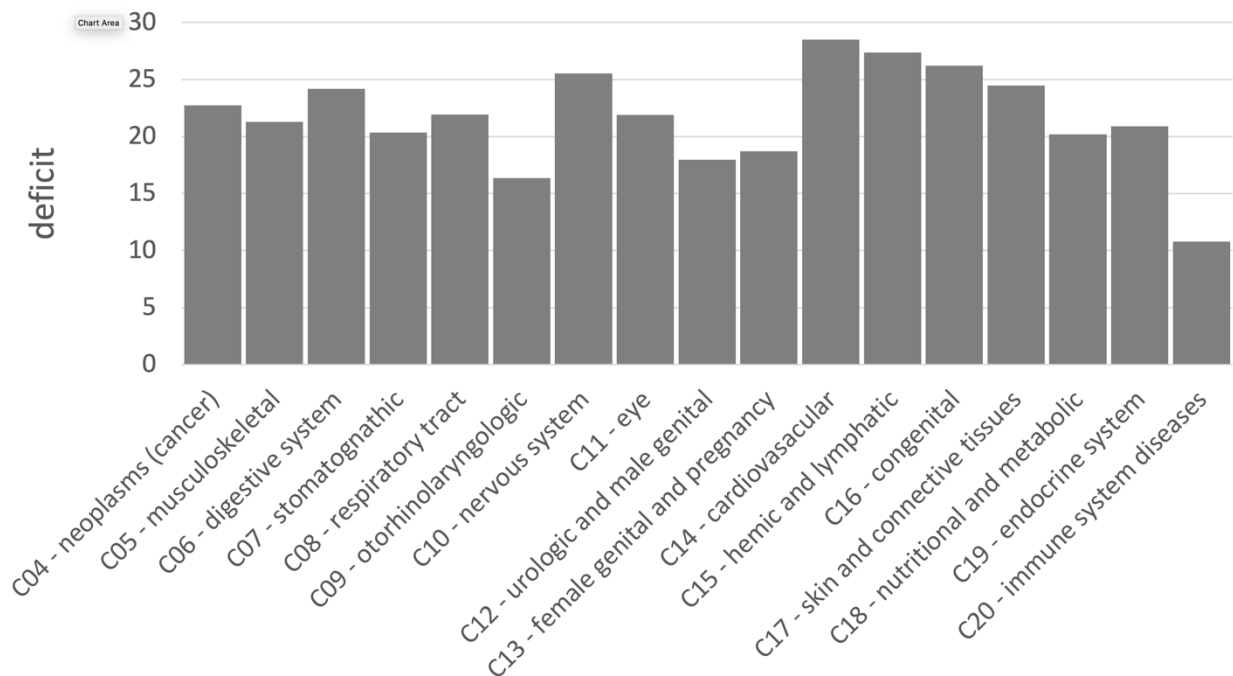**Figure 4. Sweep deficit as a function of recombination and disease variants number.**
The sweep deficit is measured as the FPR score per gene (to make all tested groups comparable) over all window sizes, and $nS_L$ and iHS, as in Figure 1 (Methods). The different groups are separated according to recombination and numbers of disease variants so that they have approximately the same size (a half or a fourth of the disease genes). All deficits are measured using only African populations.

**Similar levels of sweep depletion in disease genes across MeSH disease classes**

Because we found an overall sweep depletion at disease genes, we further ask if genes associated with different diseases might show different patterns of depletion (always in African populations). We classify disease genes into different classes according to the Medical Subject Headings (MeSH) annotation for diseases in DisGeNet (Piñero et al., 2020). The MeSH annotations organize the disease genes into 24 broad disease categories that overlap with distinct organs or large physiological systems (for example the endocrine system). We find significant (FPR<0.05) sweep depletions for all but one disease MeSH classes (FPR<0.05; Figure 5). The sweep deficit is mostly comparable across MeSH disease classes (Figure 5), suggesting that the evolutionary process at the origin of the sweep deficit is not disease-specific. This is compatible with a non-disease specific explanation such as recessive deleterious variants interfering with adaptive variants. The only non-significant deficit is for the MeSH term immune system diseases. Interestingly, there is evidence that past adaptation at disease genes in response to diverse pathogens has resulted in increased prevalence of specific auto-immune diseases (Barreiro and Quintana-Murci, 2010), and we can speculate that this is why we do not see a sweep deficit at those genes.

17

**Figure 5. Sweep deficit per MeSH disease classes.**
The sweep deficit is measured as the overall FPR score per gene (Methods), to make all MeSH classes comparable even if they include different numbers of genes.

**Discussion:**

We found a depletion of the number of genes in recent sweeps at human non-infectious, mendelian disease genes compared to non-disease genes. Although more work is now needed, the lack of sweeps at disease genes already favors specific evolutionary processes over others. For example, it makes it unlikely that past adaptations increasing the occurrence of disease variants through hitchhiking would be the dominant process linking disease and adaptation at the gene level. The lack of sweeps at disease genes also seems to be unrelated to any difference in mutation accumulation between disease and non-disease genes, since we find no sign of a difference in mutation rates between the two categories of genes in the first place, and since we match metrics accounting for mutation rate in our comparisons (for example, GC content and pS). Instead, a lack of sweeps, once selective constraint has been controlled for, seems to favor a relationship involving a lag of adaptation at disease genes beyond simple constraint (measured by the amount of deleterious mutations that are removed).

18

519     Multiple mechanisms might explain such a lag of adaptation. A first possible hypothesis is that

520     disease genes are genes that can be sensitive to the environment and whose fitness optimum can

521     change during evolution when the environment changes. However, when this happens,

522     adaptation then might take more time to chase the new optimum. Although higher pleiotropy is a

523     tempting hypothesis to explain such a lag (Otto, 2004), genes involved in multiple diseases do

524     not have a particularly pronounced sweep depletion compared to genes associated with only one

525     disease. Completely excluding pleiotropy may however require more effort, notably by

526     considering measures of pleiotropy other than the number of diseases a gene has been associated

527     with.

528

529     Another hypothesis is that disease genes may have a distribution of deleterious fitness effects

530     that is different from other genes, but that the metrics of constraint that we used do not capture

531     this difference. Specifically, we can imagine a case where disease genes have more currently

532     segregating recessive deleterious variants than other genes, and where selective sweeps are

533     impeded due to the interference of genetically linked recessive deleterious variants. The

534     deleterious effects of these variants can reveal themselves when they hitchhike together with an

535     advantageous variant that is just starting to increase in frequency (Assaf et al., 2015) .

536     Accordingly, we find a marked sweep depletion when restricting the comparison to disease and

537     non-disease genes in low recombination regions of the genome and with higher numbers of

538     disease variants (Figure 4). All these comparisons are however indirect, and we do not quantify

539     directly the amount of recessive deleterious mutations at disease or non-disease genes. Further

540     verifying that recessive deleterious mutations impede sweeps more at disease than non-disease

541     genes will require showing that recessive deleterious mutations are indeed more abundant at

542     disease genes, ideally by also estimating dominance coefficients. That said, the majority of

543     disease variants are known to be recessive and using the number of disease variants, as done in

544     the present study, should be a good proxy of the actual number of segregating recessive

545     deleterious mutations. Estimating dominance may prove challenging, since it is difficult to

546     distinguish selection coefficient changes from dominance coefficient changes (Huber et al.,

547     2018). Again, our results provide preliminary evidence to further test in the future.

548

549    In addition to suggesting possible explanatory evolutionary scenarios, our results highlight a

550    number of potential limitations and biases that also need to be explored in more detail. First, the

551    lack of sweeps at disease genes suggests the possibility of a technical bias against the annotation

552    of disease genes in sweep regions with high LD, as described in the Results. This bias is unlikely

553    to be the dominant explanation for our results, because then we would expect a stronger sweep

554    deficit at disease genes in Europe than in Africa, given that most disease genes were annotated in

555    Europe. The recombination rate at disease genes is also not different from the recombination rate

556    at non-disease genes (Figure 1). The increase of the sweep deficit when comparing disease and

557    non-disease genes only in low recombination regions (Figure 4), where disease annotation would

558    then be more difficult regardless of overlapping a sweep or not, also suggests that this bias is

559    unlikely. That said, it will still be useful to further investigate in the future how much this

560    potential bias might have contributed to our observations.

561    Second, even though more intense genetic drift seems a reasonable explanation for the less

562    pronounced sweep deficit at disease genes in Europe and East Asia than in Africa, this claim

563    needs to be further tested, for example with population simulations reproducing past population

564    demographic fluctuations. Such simulations would make it possible to test whether or not past

565    bottlenecks in ancestral Eurasian populations were strong enough to erase the sweep deficit

566    signal at disease genes in East Asia and Europe, by swamping it with random false positive

567    sweep signals.

568

569    Further work is also required regarding the connection between the sweep deficit and polygenic

570    adaptation not leaving hitchhiking signals. Our results could be explained by a general lack of

571    adaptation at disease genes, or instead by a different balance between sweeps and polygenic

572    adaptation at disease genes, with less sweeps but more polygenic adaptation that would be less

573    affected by interference with deleterious variants. It may be possible to use recent polygenic

574    adaptation quantification tools such as PALM (Stern et al., 2021) to compare its prevalence at

575    disease and non-disease genes.

576

577    Finally, there are multiple directions to further analyze the sweep deficit at disease genes that we

578    have not explored in this manuscript. For instance, analyzing the sweep deficit as a function of

579    the time of onset of diseases (early or late in life), might further provide clues to why the sweep

580    deficit exists in the first place. Preliminary comparison of the sweep deficit at specific MeSH

581    disease classes (Figure 5) with known early (congenital diseases) or mostly late onsets (cancer,

582    cardiovascular) however suggests that the average onset time of diseases might not make much

583    of a difference.

584

585    In conclusion, although our analysis reveals a strong deficit of selective sweeps at human disease

586    genes, it also suggests that more work is needed to better understand the evolutionary processes

587    at work, and the biases that may have skewed our interpretations. Despite these limitations, our

588    comparison nevertheless already suggests that specific evolutionary relationships between

589    disease genes and adaptation might be more prevalent than others, especially interference

590    between recessive deleterious and adaptive variants. As an important follow-up question, it may

591    now be important to ask how the sweep deficit at disease genes might have hidden interesting

592    adaptive patterns in previous functional enrichment analyses, especially in gene functions that

593    are often annotated based on disease evidence in the first place. For example, metabolic genes

594    are believed to be of particular interest for adaptation to climate change. But metabolic genes are

595    often found due to their role in metabolic disorders, and a strong representation of disease genes

596    among all metabolic genes could then in theory mask any sweep enrichment. A sweep

597    enrichment at metabolic genes might only become visible once controlling for the proportion of

598    disease genes, in addition to the list of controls that we already use in the present analysis

599    (Methods). Our results thus highlight the complexity of studying functional patterns of

600    adaptation in the human genome.

601

602

603

604

605

606

607

608

609

610

**Methods**

**Disease gene lists**

We consider genes that are known to be associated with diseases as disease genes. We focus on protein-coding genes associated with human mendelian non-infectious diseases. Complex diseases are associated with several loci and environmental factors. Patterns of positive selection at complex disease and mendelian disease genes may differ (Blekhman et al., 2008), which is why we restrict our analysis to mendelian disease genes. We also restrict our analyses to non-infectious disease genes, since interactions with pathogens are an entirely different problem. We nevertheless control for the proportion of genes that are immune genes or interact with viruses (see below), since it has been shown that immune genes and interactions with viruses drive a large proportion of genomic adaptation in humans (Enard et al., 2016; Castellano et al., 2019). Therefore, different proportions of immune and virus-interacting genes between disease and non-disease genes might confound their comparison. Moreover, although diseases can be associated with non-coding genes, we only use protein-coding genes. We curate disease genes defined as genes associated with diseases according to both DisGeNet (Piñero et al., 2020) and OMIM (Amberger et al., 2019), to ensure that we focus on high-confidence disease genes. DisGeNet is a comprehensive database including gene-disease associations (GDAs) from many sources. In order to get disease genes with high confidence, we further only use GDAs curated by UniProt. These gene-disease associations are extracted and carefully curated from the scientific literature and the OMIM (Online Mendelian Inheritance in Man) database, which reports phenotypes either mendelian or possibly mendelian (Amberger et al., 2019). We also exclude all genes associated with infectious diseases according to MeSH annotation (disease class C01). In the end, we curate 4215 non-infectious mendelian disease genes from DisGeNet also curated by OMIM and Uniprot. Although we rely on GDAs from Uniprot to curate high-quality disease genes, we also include GDAs of DisGeNet from other sources when classifying disease genes into different MeSH classes and measuring pleiotropy, as long as a disease gene has at least one GDA curated by OMIM and Uniprot. We completely exclude GDAs that are only reported by CTD (Comparative Toxicogenomics Database) (Davis et al., 2021) in this study. This is because CTD includes a broad range of chemical-induced diseases that might only happen where people are exposed to these chemicals, especially some inorganic chemicals that may not be present in natural environments (Davis et al., 2021).

22

642

643    In order to study different types of diseases, we also divide disease genes into different

644    classes according to the annotated MeSH classes in DisGeNet (Piñero et al., 2020). Those

645    diseases without MeSH class are annotated as "unclassfied". Genes belonging to more than one

646    MeSH class are counted in each MeSH class where they are present. MeSH classes including

647    less than 50 genes are not considered in this study. We classify all the non-infectious disease

648    genes into 22 MeSH classes including Neoplasms (C04), Musculoskeletal Diseases (C05),

649    Digestive System Diseases (C06), Stomatognathic Diseases (C07), Respiratory Tract Diseases

650    (C08), Otorhinolaryngologic Diseases (C09), Nervous System Diseases (C10), Eye Diseases

651    (C11), Male Urogenital Disease (C12), Female Urogenital Diseases and Pregnancy

652    Complications (C13), Cardiovascular Diseases (C14), Hemic and Lymphatic (C15), Congenital,

653    Hereditary, and Neonatal Diseases and Abnormalities (C16), Skin and Connective Tissue

654    Diseases (C17), Nutritional and Metabolic Diseases (C18), Endocrine System Diseases (C19),

655    Immune System Diseases (C20), Mental Disorders (F03) and "unclassified".

656

657    **Detecting selection signals at human genes**

658    All the analyses were conducted human genome version hg19. We use two different methods to

659    detect selective sweeps in human populations: iHS (integrated Haplotype Score, Voight et al.,

660    2006) and $nS_L$ (Ferrer-Admetlla et al., 2014). Both approaches are haplotype-based statistics

661    calculated with polymorphism data. We use human genome data from the 1,000 Genomes

662    Project phase 3, which includes 2,504 individuals from 26 populations (The 1000 Genomes

663    Project Consortium, 2015).

664    We measure iHS and $nS_L$ in windows centered on human coding genes (i.e. windows

665    whose center is located half-way between the most upstream transcript start site and most

666    downstream transcript stop site of protein coding genes). We use windows of sizes ranging from

667    50 kb to 1,000 kb (50kb, 100kb, 200kb, 500kb and 1,000kb) since we do not want to presuppose

668    of the size of sweeps, and since the size of the selective sweeps may vary between different

669    genes. Moreover, to avoid any preconception related to the expected strength or number of

670    sweep signals, we use a moving rank threshold strategy to measure the enrichment or deficit in

671    sweeps at disease genes. For example, we select the top 500 genes with the stronger sweep

672    signals according to a specific statistic (iHS or $nS_L$). We then compare the number of diseases

23

673    and non-disease genes within the top 500 genes with the strongest iHS or $nS_L$ signals. This was

674    repeated for different top thresholds and the corresponding ranks from top 5,000 to top 10

675    (Figure 3). Genes are ranked based on the average iHS or $nS_L$ in their gene centered windows.

676    Both iHS and $nS_L$ measure, individually for each SNP in the genome, how much larger

677    haplotypes linked to the derived SNP allele are compared to haplotypes linked to the ancestral

678    allele (Voight et al., 2006; Ferrer-Admetlla et al., 2014). For each window, we measure the

679    average of the absolute value of iHS or $nS_L$ over all the SNPs in that window with an iHS or $nS_L$

680    value. The average iHS or $nS_L$ values in a window provide high power to detect recent select

681    sweeps (Enard and Petrov, 2020).

682

683    **Comparing recent adaptation between disease and non-disease genes**

684    We use a previously developed gene-set enrichment analysis pipeline to compare recent

685    adaptation between disease and non-disease genes (Enard and Petrov, 2020)

686    (https://github.com/DavidPierreEnard/Gene_Set_Enrichment_Pipeline). This pipeline includes

687    two parts. The first part is a bootstrap test that estimates the whole sweep enrichment or

688    depletion curve at genes of interest (disease genes in our case). The second part is a false positive

689    risk (also known as false discovery rate in the context of multiple testing) that estimates the

690    statistical significance of the whole sweep enrichment curve using block-randomized genomes.

691

692    To compare disease and non-disease genes, we first need to select control non-disease genes that

693    are sufficiently far away from disease genes. In that way, we avoid using as controls non-disease

694    genes that overlap the same sweeps as neighboring disease genes, thus resulting in an

695    underpowered comparison.  The question is then how far do we need to choose non-disease

696    control genes? Ideally, we would choose non-disease control genes as far as possible from

697    disease genes in the human genome, further than the size of the largest known sweeps (for

698    example the lactase sweep), which would be on the order of a megabase. However, because there

699    are many disease genes in our dataset (4,215), there are very few non-disease genes in the human

700    genome that are more than one megabase away from the closest disease gene. This is a problem,

701    because the available number of potential control non-disease genes is an important parameter

702    that can affect both the type I error, false positive rate, and type II error, false negative rate of the

703    disease vs. non-disease genes comparison. Indeed, the smaller the control set, the more likely it

24

704     is to deviate from being representative of the true null expectation at non-disease genes. The

705     noise associated with a small sample could go either way. Either the small control sample

706     happens by chance to have less sweeps, and the bootstrap test we use to compare disease and

707     non-disease genes will become too liberal to detect sweep enrichments, and to conservative to

708     detect sweep deficits. Or the small control sample happens by chance to have more sweeps than

709     a larger control sample would, and the bootstrap test becomes too conservative to detect sweep

710     enrichments, and too liberal to detect sweep deficits.

711     After trying distances between disease genes and control disease genes of 100kb, 200kb, 300kb,

712     400kb and 500kb, we find that the sweep deficit observed at disease genes increases steadily

713     from 100kb to 300kb (Table 1), showing that 100kb or 200kb are likely insufficient distances.

714     Further than 300kb at 400kb, we do not observe much stronger sweep deficits than at 300kb,

715     while at the same time the risks of type I and type II errors keep increasing due to shrinking non-

716     disease genes control sets. This would translate in a decreased power to possibly exclude the null

717     hypothesis of no sweep enrichment or deficit in the second part of the pipeline, when estimating

718     the actual pipeline FPR. Because of this, we set the required distance of potential control non-

719     disease genes from disease genes at 300kb. This is also the distance where there are still

720     approximately as many control genes (3455) as there are disease genes that we can use for the

721     comparison (3030; those genes out of the 4,215 disease genes with sweep data and data for all

722     the confounding factors).

723

724

| minimal distance | sweep deficit |
|---|---|
| 100kb | -20889 |
| 200kb | -35009 |
| 300kb | -68928 |
| 400kb | -88546 |

725     **Table 1. Sweep deficit as a function of the minimal distance of control non-disease genes.**
726     The sweep deficit is measured by the FPR score, that is the cumulative difference between the
727     number of genes in sweeps at disease and control non-disease genes, across window sizes, sweep
728     summary statistics, and African populations (see the rest of the Methods).

729

730

731    Another important aspect of the bootstrap test (first part of the pipeline), aside from setting up

732    the minimal distance of the control non-disease genes, is the matching of potential confounding

733    factors likely to influence sweep occurrence. We choose non-disease control genes that have the

734    same confounding factors characteristics as disease genes (for example, control non-disease

735    genes that have the same gene expression level across tissues as disease genes). The precise

736    matching algorithm is detailed in Enard & Petrov (2020).

737    When comparing disease and non-disease genes with the bootstrap test, we control for the

738    following potential confounding factors that could influence the occurrence of sweeps at genes:

739    ● Average overall expression in 53 GTEx v7 tissues (The GTEx Consortium, 2015)

740        (https://www.gtexportal.org/home/). We used the log (in base 2) of TPM (Transcripts Per

741        Million).

742    ● Expression (log base 2 of TPM) in GTEx lymphocytes.  Expression in immune tissues

743        may impact the rate of sweeps.

744    ● Expression (log base 2 of TPM) in GTEx testis. Expression in testis might also impact the

745        rate of sweeps.

746    ● deCode recombination rates 50kb and 500kb: recombination is expected to have a strong

747        impact on iHS and $nS_L$ values, with larger, easier to detect sweeps in low recombination

748        regions but also more false positive sweeps signals. The average recombination rates in

749        the gene-centered windows are calculated using the most recent deCode recombination

750        map (Halldorsson et al., 2019). We use both 50kb and 500kb window estimates to

751        account for the effect of varying window sizes on the estimation of this confounding

752        factor (same logic for other factors where we also use both 50kb and 500kb windows).

753    ● GC content is calculated as a percentage per window in 50kb and 500kb windows. It is

754        obtained from the USCS Genome Browser (Kent et al., 2002).

755    ● The density of coding sequences in 50kb and 500kb windows centered on genes. The

756        density is calculated as the proportion of coding bases respect to the whole length of the

757        window. Coding sequences are Ensembl v99 coding sequences.

758    ● The density of mammalian phastCons conserved elements (Siepel et al., 2005) (in 50kb

759        and 500k windows), downloaded from the UCSC Genome Browser (Kent et al., 2002).

760        We used a threshold considering 10% of genome as conserved, as it is unlikely that more

761        than 10% of the whole genome is constrained according to previous evidence (Siepel et

26

762      al., 2005). Given that each conserved segment had a score, we considered those segments
763      above the 10% threshold as conserved.

764   ● The density of regulatory elements, as measured by the density of DNASE1
765      hypersensitive sites (in 50kb and 500kb windows) also from the UCSC Genome Browser
766      (Kent et al., 2002).

767   ● The number of protein-protein interactions (PPIs) in the human protein interaction
768      network (Luisi et al., 2015). The number of PPIs has been shown to influence the rate of
769      sweeps (Luisi et al., 2015). We use the log (base 2) of the number of PPIs.

770   ● The gene genomic length, i.e. the distance between the most upstream and the most
771      downstream transcription start sites.

772   ● The number of gene neighbors in a 50kb window, and the same number in 500kb window
773      centered on the focal genes: it is the number of coding genes within 25kb or within
774      250kb.

775   ● The number of viruses that interact with a specific gene (Enard and Petrov, 2020).

776   ● The proportion of immune genes. The matched control sets have the same proportion of
777      immune genes as disease genes, immune genes being genes annotated with the Gene
778      Ontology terms GO:0002376 (immune system process), GO:0006952 (defense response)
779      and/or GO:0006955 (immune response) as of May 2020 (Gene Ontology Consortium,
780      2021).

781   ● The average number of non-synonymous variants PN in African populations, and the
782      number of synonymous variants PS. We matched PN to build control sets of non-disease
783      genes with the same average amount of strong purifying selection as disease genes. Also,
784      PS can be a proxy for mutation rate and we can build control sets of non-disease genes
785      with similar level of mutation rates.

786   ● McVicker's B value which can be used to account for the effect of background selection
787      on rates of adaptation and especially weak adaptation (McVicker et al., 2009).

788

789 Similar to the selection of control genes far enough from disease genes, the matching of many
790 confounding factors decreases the number of non-disease genes that can effectively be used as
791 controls. This further increases the risk of type I and type II errors of the bootstrap test, as
792 previously described. In addition, the bootstrap test only provides p-value for each tested sweep

27

793    rank threshold separately, in the whole enrichment (or deficit) curve (Figure 2). It does not

794    provide any estimate of the significance of the whole curve, which is needed to estimate the

795    significance of a sweep enrichment or deficit without making too many assumptions on how

796    many sweeps are expected or how strong they are.

797    To address the increased type I and type II error risks of the bootstrap test, as well to get an

798    unbiased significance estimate for whole enrichment curves, the second part of our pipeline

799    conducts a false positive risk analysis based on block-randomized genomes (Enard and Petrov,

800    2020). Briefly, we re-estimate many whole enrichment curves reusing the same disease and

801    control non-disease genes used in the first part of the pipeline by the bootstrap test, but after

802    having randomly shuffled the locations of genes or clusters of neighboring genes in sweeps at

803    those disease and control non-disease genes. To do this, we order the disease and control non-

804    disease genes as they appear in the genome. We then define blocks of neighboring genes, whose

805    limits do not interrupt clusters of genes in the same putative sweep. Then, we randomly shuffle

806    the order of these blocks. Because we do not cut any cluster of genes that might be in the same

807    sweep, the resulting block-randomized genomes preserve the same clustering of the genes in the

808    same putative sweeps as in the real genome. With this approach, we look at the exact same set of

809    disease and control non-disease genes and just shuffle sweep locations between them. Thus, by

810    using many block-randomized genomes, we can estimate the null expected range of whole

811    enrichment curves while fully accounting for the extra variance expected from having a limited

812    sample of control non-disease genes. We can then estimate a false positive risk (FPR) for the

813    whole enrichment or deficit curve by comparing the real observed one with the distribution of

814    random curves generated with block-randomized genomes.

815

816    To measure the FPR for a curve, we need to define a metric to compare the real curve with the

817    randomly generated ones. In figure 1, we show relative enrichments at each sweep rank

818    threshold, the number of disease genes in sweeps divided by the number of control non-disease

819    genes in sweeps. As a summary metric for the curve, we could then use the sum of the relative

820    enrichments over all thresholds. However, the issue with this approach is that a relative

821    enrichment is the same whether we have 2 disease genes in sweeps and one control non-disease

822    gene in sweeps, or we have 200 disease genes in sweeps and 100 control non-disease genes in

823    sweeps. Thus, although relative enrichments are convenient for visualization on a figure, they are

824    not adequate to measure the FPR. Instead of the relative enrichment, we use the difference

825    between disease and non-disease genes, that is, the number of disease genes in sweeps, minus the

826    average number of control non-disease genes across control sets built by the bootstrap test. We

827    then use as a metric for a whole curve the sum of differences over all the rank thresholds. We use

828    this sum of differences to estimate the enrichment or deficit curve FPR, as the proportion of

829    block-randomized genomes where the sum of differences exceeds the observed sum of

830    differences for an enrichment (one minus this proportion for a deficit).

831

832    Importantly, although so far we have described the case where we measure the FPR for one

833    enrichment curve, nothing prevents us from calculating a single sum of differences over an entire

834    group of enrichment or deficit curves. This way, we can measure a single FPR for any number of

835    curves considered together. In our analysis, we measure a single FPR adding iHS and $nS_L$ curves

836    together, and also adding together the curves for 50kb, 100kb, 200kb, 500kb and 1000kb

837    windows (ten curves in total, 2 statistics*5 window sizes).

838

839    **Sweep deficit at high and low recombination disease genes, and at high and low disease**

840    **variant number disease genes**

841    To generate Figure 4, we separate disease genes in groups of approximately the same size based

842    on their recombination rate and numbers of disease variants annotated in OMIM/Uniprot. We

843    separate the disease genes into two groups of equal size, those with recombination lower than

844    1.137 cM/Mb, and those with recombination higher than this value. To count the disease variants

845    at each disease gene, we count not only the OMIM/Uniprot disease variants for that gene, but

846    also all the other OMIM/Uniprot disease variants that occur in a 500kb window centered on that

847    gene. We do this because the recessive deleterious variants form other nearby disease genes may

848    also interfere with adaptation. Half of disease genes have less than five OMIM/Uniprot disease

849    variants, and half have five or more.

850

851    **Impact of functional differences between disease and non-disease genes on the sweep deficit**

852    The sweep deficit at disease genes could be due to a different representation of gene functions at

853    disease genes compared to control non-disease genes. In this case, disease genes would have less

854    adaptation not because they are disease genes, but because the gene functions that are enriched

29

855    among disease genes compared to non-disease happen to experience less adaptation. We can test

856    this possibility using Gene Ontology (GO) (Gene Ontology Consortium, 2021) functional

857    annotations as follows. If GO gene functions that are enriched in disease genes experience less

858    adaptation independently of the disease status of genes, then we can predict that non-disease

859    genes with these functions should also experience less adaptation than non-disease genes that do

860    not have these GO functions. In total, we find that 3,097 GO annotations are enriched in disease

861    genes compared to confounding factors-matched controls (bootstrap test P≤0.01). In our dataset,

862    half of non-disease genes have 20 or more of these GO annotations, and half have less than

863    twenty (very few have none). We find no difference in the sweep prevalence between the two

864    groups (20 or more annotations vs. less than 20 annotations at least 300kb away; FPR=0.15). The

865    sweep deficit at disease genes is therefore unlikely to be due to the gene functions that are more

866    represented in disease genes compared to controls. In addition, such a scenario would not explain

867    the lack of sweep deficit observed at disease genes with high recombination rates and low

868    numbers of disease variants (Figure 4).

869

**Acknowledgements**

871    We wish to thank Dan Shrider for helpful comments on the results presented in the manuscript.

872

**Author Contributions**

874    Conceived and designed the analyses: CD, DE. Performed the analyses: CD and DE. Wrote the

875    manuscript: CD, DST and DE. Interpreted the results: CD, DST, MEL and DE.

876

**References**

878

879    Amberger JS, Bocchini CA, Scott AF, Hamosh A. 2019. OMIM.org: leveraging knowledge

880        across phenotype-gene relationships. *Nucleic Acids Res*. 47(D1):D1038–D1043.

881    Assaf ZJ, Petrov DA, Blundell JR. 2015. Obstruction of adaptation in diploids by recessive,

882        strongly deleterious alleles. *Proc. Natl. Acad. Sci. U. S. A.* 112 (20):E2658-E2666.

883    Barreiro LB, Quintana-Murci L. 2010. From evolutionary genetics to human immunology: how

884        selection shapes host defence genes. *Nat. Rev. Genet.* 11(1):17–30.

885    Birky CW, Walsh JB. 1988. Effects of linkage on rates of molecular evolution. *Proc. Natl. Acad.*

886        *Sci. U. S. A.* 85:6414–6418.

887   Blekhman R, Man O, Herrmann L, Boyko AR, Indap A, Kosiol C, Bustamante CD, Teshima
888         KM, Przeworski M. 2008. Natural Selection on Genes that Underlie Human Disease
889         Susceptibility. *Curr. Biol.* 18(12):883–889.

890   Castellano D, Uricchio LH, Munch K, Enard D. 2019. Viruses rule over adaptation in conserved
891         human proteins. *bioRxiv*, 555060.

892   Chun S, Fay JC. 2011. Evidence for Hitchhiking of Deleterious Mutations within the Human
893         Genome. *PLOS Genet.* 7(8): e1002240.

894   Davis AP, Grondin CJ, Johnson RJ, Sciaky D, Wiegers J, Wiegers TC, Mattingly CJ. 2021.
895         Comparative Toxicogenomics Database (CTD): update 2021. *Nucleic Acids Res.*
896         49(D1):D1138–D1143.

897   Enard D, Cai L, Gwennap C, Petrov DA. 2016. Viruses are a dominant driver of protein
898         adaptation in mammals. *eLlife* 5: e12469.

899   Enard D, Messer PW, Petrov DA. 2014. Genome-wide signals of positive selection in human
900         evolution. *Genome Res.* 24(6):885-95.

901   Enard D, Petrov DA. 2020. Ancient RNA virus epidemics through the lens of recent adaptation
902         in human genomes. *Philos. Trans. R. Soc. B Biol. Sci.* 375:20190575.

903   Eyre-Walker YC, Eyre-Walker A. 2014. The Role of Mutation Rate Variation and Genetic
904         Diversity in the Architecture of Human Disease. *PLoS One* 9(2):e90166.

905   Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R. 2014. On detecting incomplete soft or
906         hard selective sweeps using haplotype structure. *Mol. Biol. Evol.* 31(5):1275–1291.

907   Gene Ontology Consortium. 2021. The Gene Ontology resource: enriching a GOld mine.
908         *Nucleic Acids Res.* 49(D1):D325–D334.

909   Halldorsson B V, Palsson G, Stefansson OA, Jonsson H, Hardarson MT, Eggertsson HP,
910         Gunnarsson B, Oddsson A, Halldorsson GH, Zink F, et al. 2019. Characterizing mutagenic
911         effects of recombination through a sequence-level genetic map. *Science*
912         363(6425):eaau1043.

913   Hill WG, Robertson A. 1966. The effect of linkage on limits to artificial selection. *Genet. Res.*
914         8(3):269–294.

915   Huber CD, Durvasula A, Hancock AM, Lohmueller KE. 2018. Gene expression drives the
916         evolution of dominance. *Nat. Commun.* 9:2750.

917   Ittisoponpisan S, Alhuzimi E, Sternberg MJE, David A. 2017. Landscape of Pleiotropic Proteins

918     Causing Human Disease: Structural and System Biology Insights. *Hum. Mutat.* 38(3):289–
919     296.

920   Jain K. 2019. Interference Effects of Deleterious and Beneficial Mutations in Large Asexual
921     Populations. *Genetics* 211(4):1357–1369.

922   Johnson T, Barton NH. 2002. The effect of deleterious alleles on adaptation in asexual
923     populations. *Genetics* 162(1):395–411.

924   Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler D. 2002. The
925     human genome browser at UCSC. *Genome Res.* 12(6):996–1006.

926   Kim PM, Korbel JO, Gerstein MB. 2007. Positive selection at the protein network periphery:
927     Evaluation in terms of structural constraints and cellular context. *Proc. Natl. Acad. Sci. U.
928     S. A.* 104(51):20274–20279.

929   Luisi, P., Alvarez-Ponce, D., Pybus, M., Fares, M. A., Bertranpetit, J., & Laayouni, H. 2015.
930     Recent positive selection has acted on genes encoding proteins with more interactions
931     within the whole human interactome. *Genome Bio. Evol.* 7(4):1141–1154.

932   McVicker G, Gordon D, Davis C, Green P. 2009. Widespread genomic signatures of natural
933     selection in hominid evolution. *PLoS Genet.* 5: e1000471.

934   Osada N, Mano S, Gojobori J. 2009. Quantifying dominance and deleterious effect on human
935     disease genes. *Proc. Natl. Acad. Sci. U. S. A.* 106(3):841 – 846.

936   Otto SP. 2004. Two steps forward, one step back: the pleiotropic effects of favoured alleles.
937     *Proc. Biol. Sci.* 271(1540):705–714.

938   Park S, Yang JS, Kim J, Shin YE, Hwang J, Park J, Jang SK, Kim S. 2012. Evolutionary history
939     of human disease genes reveals phenotypic connections and comorbidity among genetic
940     diseases. *Sci. Rep.* 2:757.

941   Peck JR. 1994. A ruby in the rubbish: beneficial mutations, deleterious mutations and the
942     evolution of sex. *Genetics* 137(2):597–606.

943   Piñero J, Ramírez-Anguita JM, Saüch-Pitarch J, Ronzano F, Centeno E, Sanz F, Furlong LI.
944     2020. The DisGeNET knowledge platform for disease genomics: 2019 update. *Nucleic
945     Acids Res.* 48(D1):D845–D855.

946   Quintana-murci L. 2016. Understanding rare and common diseases in the context of human
947     evolution. *Genome Biol.*17:225.

948   Schrider DR. 2020. Background Selection Does Not Mimic the Patterns of Genetic Diversity

949       Produced by Selective Sweeps. *Genetics* 216(2): 499-519.

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

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

952       W. J., Miller, W., & Haussler, D. 2005. Evolutionarily conserved elements in vertebrate,

953       insect, worm, and yeast genomes. *Genome Res*. 15(8):1034–1050.

954 Smith JM. 1976. What Determines the Rate of Evolution? *Am. Nat*. 110(973):331–338.

955 Smith JM, Haigh J. 1974. The hitch-hiking effect of a favourable gene. *Genet. Res*. 23:23–35.

956 Smith NGC, Eyre-Walker A. 2003. What Determines the Rate of Evolution? *Gene* 318:169–175.

957 Spataro N, Rodríguez JA, Navarro A, Bosch E. 2017. Properties of human disease genes and the

958       role of genes linked to Mendelian disorders in complex disease aetiology. *Hum. Mol. Genet*.

959       26(3):489–500.

960 Stern AJ, Speidel L, Zaitlen NA, Nielsen R. 2021. Disentangling selection on genetically

961       correlated polygenic traits via whole-genome genealogies. *Am. J. Hum. Genet*. 108(2):219–

962       239.

963 The 1000 Genomes Project Consortium. 2015. A global reference for human genetic variation.

964       *Nature* 526:68–74.

965 The GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue

966       gene regulation in humans. 2015. *Science*. 348(6235):648-660

967 Uricchio LH, Petrov DA, Enard D. 2019. Exploiting selection at linked sites to infer the rate and

968       strength of adaptation. *Nat. Ecol. Evol*. 3:977–984.

969 Voight BF, Kudaravalli S, Wen X, Pritchard JK. 2006. A map of recent positive selection in the

970       human genome. *PLoS Biol*. 4(3):e72.

971 Xu J, Zhang J. 2014. Why Human Disease-Associated Residues Appear as the Wild-Type in

972       Other Species: Genome-Scale Structural Evidence for the Compensation Hypothesis. *Mol.*

973       *Biol. Evol*. 31(7):1787–1792.