

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

Excess of Deleterious Mutations around HLA Genes Reveals Evolutionary Cost of Balancing Selection

Tobias L. Lenz^{*,1,2}, Victor Spirin^{†,1}, Daniel M. Jordan^{‡,1} & Shamil R. Sunyaev^{*,1,3}

Affiliations:

¹ Division of Genetics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA 02115, USA

² Evolutionary Immunogenomics, Department of Evolutionary Ecology, Max Planck Institute for Evolutionary Biology, 24306 Plön, Germany

³ Program in Medical and Population Genetics, The Broad Institute, Cambridge, MA 02142, USA

[†] Present address: Merck Research Laboratories, 470 Atlantic Avenue, Boston, MA 02210, USA

[‡] Present address: Icahn Institute for Genomics and Multiscale Biology, Department of Genetics and Genomics, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA

* Corresponding authors:

Tobias L. Lenz, Email: lenz@post.harvard.edu

Shamil R. Sunyaev, Email: ssunyaev@rics.bwh.harvard.edu

Keywords: balancing selection; exome; simulations; deleterious variation; mutation load; MHC/HLA

29 **Abstract**

30 Deleterious mutations are expected to evolve under negative selection and are usually
31 purged from the population. However, deleterious alleles segregate in the human
32 population and some disease-associated variants are maintained at considerable
33 frequencies. Here we test the hypothesis that balancing selection may counteract
34 purifying selection in neighboring regions and thus maintain deleterious variants at
35 higher frequency than expected from their detrimental fitness effect. We first show in
36 realistic simulations that balancing selection reduces the density of polymorphic sites
37 surrounding a locus under balancing selection, but at the same time markedly increases
38 the population frequency of the remaining variants, including even substantially
39 deleterious alleles. To test the predictions of our simulations empirically, we then use
40 whole exome sequencing data from 6,500 human individuals and focus on the most
41 established example for balancing selection in the human genome, the major
42 histocompatibility complex (MHC). Our analysis shows an elevated frequency of
43 putatively deleterious coding variants in non-HLA genes localized in the MHC region.
44 The mean frequency of these variants declined with physical distance from the classical
45 HLA genes, indicating dependency on genetic linkage. These results reveal an indirect
46 cost of the genetic diversity maintained by balancing selection, which has hitherto been
47 perceived as mostly advantageous, and have implications both for the evolution of
48 recombination and also for the epidemiology of various MHC-associated diseases.

49
50

51 **Introduction**

52 A large number of population genetics studies point to the existence of numerous mildly
53 deleterious alleles segregating in the human population (Henn et al. 2015). This
54 abundance of deleterious alleles is evident from the comparison of human DNA
55 polymorphism to human-chimpanzee sequence divergence (Bustamante et al. 2005),
56 from the analysis of allele frequency distribution (Kryukov et al. 2007; Boyko et al. 2008;
57 Kryukov et al. 2009) and estimated allelic ages (Kiezun et al. 2013). Deleterious alleles
58 are present in individual human genomes at functionally significant sites in both protein
59 coding genes and regulatory regions (Maurano et al. 2012; Fu et al. 2013).
60 Understanding the maintenance of deleterious variation in the population is critically
61 important for evolutionary models of complex traits including common human diseases.
62 Persistence of common diseases is paradoxical from the evolutionary standpoint, and
63 the stable existence of deleterious variation in spite of the action of purifying selection
64 requires an explanation.

65
66 Different hypotheses have been put forward to explain the occurrence of deleterious
67 genetic variants in the population. Mutation-selection balance in combination with
68 demography and genetic drift certainly account for a significant part of the genetic load
69 and for the persistence of rare deleterious variation (Dudley et al. 2012; Fu et al. 2013).
70 The occurrence of more frequent variants, on the other hand, is commonly hypothesized
71 to involve selective forces favoring genetic diversity. However, selection should not
72 necessarily target the variants in question. With the improved understanding of the
73 recombination landscape and of the extent of genetic linkage, it is becoming increasingly
74 clear that many sites in the genome do not evolve completely independently and allele
75 frequency changes may be influenced by variation at neighboring linked sites (Gillespie
76 2000; Charlesworth et al. 2003; Barton 2010).

77
78 Previous research has already addressed the evolutionary fate of genetic variation
79 around sites under selection, focusing mostly on neutral variation. Background selection
80 has been shown to result in the reduction of neutral diversity surrounding regions under
81 purifying selection (Charlesworth et al. 1993). Genetic hitchhiking events associated with

82 positive selection generally eliminate surrounding genetic variation. On the other hand,
83 they transiently increase allele frequencies of variants residing on the same haplotype
84 as the positively selected variant (Maynard Smith and Haigh 1974; Fay and Wu 2000;
85 Chun and Fay 2011), and theoretical analysis and empirical data suggest that
86 hitchhiking events can also elevate the frequency of deleterious variants (Chun and Fay
87 2011; Hartfield and Otto 2011; Marsden et al. 2016). Conversely, existing recessive
88 deleterious variation may also slow the fixation of beneficial alleles on the same
89 haplotype (Assaf et al. 2015). However, since hitchhiking eventually reduces the
90 absolute amount of genetic diversity, it is not sufficient to explain regions with generally
91 elevated levels of sequence diversity across the genome.

92
93 In contrast, balancing selection leads to a long-term persistence of common genetic
94 variation in surrounding loci (Charlesworth 2006; Gao et al. 2015). Consequently, it has
95 been proposed that balancing selection may also lead to an excess of deleterious
96 variation around the balanced locus. This scenario is supported by studies of neutral
97 variation that showed increased variation in regions linked to loci under both simple and
98 multi-locus balancing selection (Kaplan et al. 1988; Grimsley et al. 1998; O'hUigin et al.
99 2000; Navarro and Barton 2002). However, a comprehensive theoretical and empirical
100 investigation of this scenario for deleterious variants, which is of critical interest for both
101 population genetics and medical genetics, is lacking so far.

102
103 The prevalence of balancing selection in the human population and its role in shaping
104 genetic variation is still debated. In contrast to models developed in the 1960s (e.g
105 Lewontin and Kojima 1960), studies of the genomic era have generally assumed that
106 balancing selection is an exception rather than the rule (Asthana et al. 2005; Bubb et al.
107 2006). However, an accumulating number of recent studies are identifying genomic
108 features characteristic of balancing selection (Bustamante et al. 2005; Andrés et al.
109 2009; Fumagalli et al. 2009; Sellis et al. 2011; Gokcumen et al. 2013; Leffler et al. 2013;
110 DeGiorgio et al. 2014; Teixeira et al. 2015).

111

112 A classical locus to investigate the effect of balancing selection on multiple sites and its
113 influence on neighboring variation is the major histocompatibility complex (MHC). The
114 MHC is one of the most prominent examples of balancing selection in the vertebrate
115 genome. This is thought to be due to the so-called classical MHC genes (in humans
116 called Human Leukocyte Antigen; HLA), whose products present antigenic peptides on
117 the cell surface and by that play a key role in the adaptive immune response (Neefjes et
118 al. 2011). These classical HLA genes are scattered across the MHC region and exhibit
119 exceptional allelic polymorphism and an extreme level of heterozygosity, which is
120 thought to increase pathogen resistance and thus to be maintained by pathogen-
121 mediated balancing selection (Trowsdale 2011). Based on well-documented signatures
122 of balancing selection at the classical HLA genes, this gene complex thus provides a
123 perfect study system to explore the effect of balancing selection on the frequency of
124 linked deleterious variation. In fact, a striking feature of the MHC, potentially related to
125 the maintenance of deleterious variation, is that the MHC region is highly enriched in
126 variants associated with common human diseases identified by Genome-Wide
127 Association Studies (GWAS; **fig. 1**). Many of the GWAS peaks are not caused by
128 variation in classical HLA genes, and a significant number of phenotypes associated
129 with genetic variants within the MHC region are not classic autoimmune or infectious
130 diseases (Trowsdale 2011).

131
132 To investigate the potential of balancing selection for maintaining deleterious genetic
133 variation, we first use a forward simulation approach and then analyze empirical data to
134 test the predictions of the simulations. Previously, the large scale sequencing data
135 required to compare the mutational load in the MHC region against the rest of the
136 genome has been lacking. Here we make use of a whole exome sequence data set of
137 about 6,500 individuals, which provides detailed population-level information of coding
138 sequence variation across the entire human genome.

139

140

141

142

143 **Results**

144 **Simulating multi-locus balancing selection**

145 We used forward simulations to investigate the effect of multi-locus balancing selection
146 at HLA genes on the evolution of deleterious genetic variants throughout the MHC
147 region. The simulations were based on empirical parameter values from the human
148 MHC region and included a central HLA gene, surrounded by non-HLA regions that
149 represented the variation in neighboring non-HLA genes. In order to simulate multi-locus
150 selection, multiple scattered sites along the virtual HLA gene evolved under different
151 HLA selection scenarios: either neutrally, under balancing selection, or, as a contrast,
152 under recurrent sweeps of positive selection. We then explored the effect of these
153 different HLA selection scenarios on the evolution of variants in HLA-neighboring
154 regions that evolved themselves under different levels of purifying selection.

155
156 Our simulations showed that multi-locus balancing selection at HLA generally reduces
157 the number of segregating sites in the neighboring sequence regions compared to a
158 scenario with no selection at HLA (**fig. 2a**). With moderate or high levels of purifying
159 selection acting at neighboring sites, the number of segregating sites was even lower.
160 Interestingly, the reduction in the density of segregating sites by balancing selection was
161 mainly due to the removal of sites with rare variants (**fig. 2b**). And as nucleotide diversity
162 (π) is strongly dependent on variant frequencies, the removal of rare variants had little
163 effect on overall diversity in the neighboring regions. Instead, balancing selection at the
164 HLA led to a significant increase in nucleotide diversity in the neighboring regions
165 surrounding the HLA and was elevated even in the presence of moderate purifying
166 selection (**fig. 2c**). Nucleotide diversity increased in spite of the reduction in number of
167 segregating sites because derived allele frequencies at the remaining sites were, on
168 average, strongly elevated, resulting in a site frequency spectrum that was enriched with
169 intermediate frequency alleles (**fig. 3a**). This increase in the number of segregating sites
170 with intermediate frequency alleles then caused an overall increase in nucleotide
171 diversity, irrespective of the significant loss of rare variants around the HLA gene. As
172 expected, increasing purifying selection in the surrounding region reduced allele
173 frequencies overall. However, relative to the other HLA selection scenarios, balancing

174 selection at the HLA in all simulated cases maintained a significant fraction of
175 deleterious variation at moderate frequencies counteracting the effect of purifying
176 selection (Wilcoxon rank-sum test, all $P < 0.001$; **fig. 3b-d**). The effect of balancing
177 selection on variant frequency in the neighboring regions was strongly dependent on the
178 rate of recombination: Artificially increasing the recombination rate by an order of
179 magnitude basically removed the effect, while an equivalent reduction in recombination
180 led to a substantially stronger enrichment for intermediate frequency variants at all levels
181 of purifying selection (**supplementary figs. S1 & S2**, Supplementary Material online).

182
183 We contrasted these results with the alternative scenario of recurrent complete selective
184 sweeps at the HLA due to continuous environmental change, e.g. constantly adapting or
185 newly arising pathogens with ‘sweeping’ selective effects (e.g. the Bubonic plague).
186 Similarly to the balancing selection scenario, recurrent sweeps in HLA genes led to a
187 reduction in the number of segregating sites in the neighboring regions (**fig. 2a**).
188 Compared to the balancing selection scenario, this reduction in SNP density is stronger
189 if there is no purifying selection on the surrounding regions and weaker when purifying
190 selection exerts pressure on the neighboring sequence. This reduction is primarily
191 affecting SNPs with common derived alleles, with almost no removal of rare variants
192 (**fig. 2b**). Consequently, in sharp contrast with the balancing selection scenario,
193 nucleotide diversity is reduced rather than elevated (**fig. 2c**).

194
195 Overall, our simulations showed that multi-locus balancing selection can lead to a
196 distinct signature of a reduced number of polymorphic sites with elevated allele
197 frequencies. The elevated site frequency spectrum suggests that balancing selection
198 has the realistic potential to significantly increase of even frequencies of strongly
199 deleterious variants in regions around the classical HLA genes.

200
201 **Exome sequencing data**

202 We used SNP variation data from the whole-exome sequence dataset of the NHLBI GO
203 Exome Sequencing Project (from a total of about 6,500 humans and 17,684 genes,
204 including 124 genes within the MHC region). Excluding eleven of the latter that are

205 classical HLA genes, expected to evolve under balancing selection, our analyses of the
206 MHC region included variant data from 113 genes. The median number of polymorphic
207 sites averaged 48 for these genes, significantly lower than the exome-wide median of 74
208 (Wilcoxon rank-sum test, $P < 0.001$; **fig. 4a**) and was mainly due to a difference in the
209 number of sites with low frequency variants (**fig. 4b**). Our simulations also showed fewer
210 polymorphic sites under the balancing selection scenarios than without it. The reduced
211 number of polymorphic sites affected all SNP categories, including variants likely to be
212 deleterious (as predicted by PolyPhen-2 (Adzhubey et al. 2010); this distinguishes the
213 effect of balancing selection from other scenarios, such as genetic hitchhiking during
214 selective sweeps where the higher frequency synonymous variants are preferentially
215 removed from the population (Chun and Fay 2011). In contrast, the population frequency
216 of derived alleles at polymorphic sites was substantially elevated within the MHC region
217 compared to the rest of the genome (**fig. 5a**).

218
219 This result is not caused by local differences in mutation rate, since the rate of *de novo*
220 mutations (as estimated from trio sequencing; Francioli et al. 2015) in the MHC region
221 does not differ from the average exome-wide rate (Wilcoxon rank-sum test, MHC vs.
222 whole exome, $P > 0.05$). The frequencies of derived alleles correlated with their position
223 relative to the nearest classical HLA locus: derived alleles at sites physically close to an
224 HLA locus were at higher frequencies than at sites further away (Spearman's $\rho = -0.14$,
225 $P < 0.001$; **fig. 5b**), suggesting that their frequencies are influenced by linkage
226 disequilibrium to HLA genes.

227
228 Since we were especially interested in potentially deleterious variants, we focused our
229 analyses on loss-of-function variants and missense variants classified as “probably
230 damaging” by PolyPhen-2 (Adzhubey et al. 2010), which are expected to evolve under
231 moderate to strong purifying selection. This expectation was supported by a
232 substantially lower exome-wide derived allele frequency of “probably damaging”
233 variants, compared to variants classified as “benign” (Wilcoxon rank-sum test, both $P <$
234 0.001 ; **supplementary fig. S3**, Supplementary Material online). We found a significant
235 elevation in average derived allele frequencies within the MHC region, for both probably

236 damaging missense variants and loss-of-function variants, compared to the average
237 frequency across the entire exome (Wilcoxon rank-sum test, both $P < 0.001$; **fig. 6**). This
238 shift in the variant frequencies was observed within both African American and
239 European American subpopulations of the ESP6500 data (**supplementary note S1**,
240 **supplementary table S1 and supplementary fig. S4**, Supplementary Material online).

241
242 Overall, we identified a substantial number of non-HLA genes throughout the MHC
243 region, carrying deleterious variants at high frequencies (**supplementary table S2**,
244 Supplementary Material online). For some of these genes, the average frequency of
245 probably damaging variants was more than two orders of magnitude above the genome-
246 wide average. Such genes include *MICA* (6 variants predicted to be ‘probably damaging’
247 with mean derived allele frequency, or DAF, of 9.7%), *PSORS1C1* (2 variants with DAF
248 of 8.8%), and *CFB* (9 variants with DAF of 2.0%). **Supplementary table S2**
249 (Supplementary Material online) lists potential diseases reported in connection with
250 these genes, including common autoimmune diseases (e.g. psoriasis, rheumatoid
251 arthritis), cancer (e.g. prostate cancer, hepatocellular carcinoma), and mental disorders
252 (e.g. Alzheimer’s disease, schizophrenia).

253

254

255 **Discussion**

256 Our simulations of genomic regions around a classical HLA gene showed that balancing
257 selection has a similar effect on the number of linked variants as recurrent sweeps of
258 positive selection in the way it reduces the absolute number of segregating sites.
259 However, while positive selection simultaneously reduces the frequency of variants at
260 the remaining sites, balancing selection has the opposite effect and substantially
261 increases derived allele frequencies, leading to a site frequency spectrum with an
262 excess of intermediate frequency alleles. As expected, derived allele frequencies also
263 depended on the strength of selection against the deleterious variants. This is in
264 agreement with previous studies on the evolution of neutral variation around balanced
265 loci, showing generally elevated sequence diversity (Kaplan et al. 1988; Grimsley et al.

266 1998; Horton et al. 1998; O'hUigin et al. 2000; Navarro and Barton 2002; Connallon and
267 Clark 2013).

268
269 Results from deep population sequencing data supported our simulation results
270 empirically, by showing the same pattern of a reduced number of polymorphic sites but
271 elevated allele frequencies in the MHC region compared to the rest of the exome. This
272 observation held for both loss-of-function mutations, which are likely to be highly
273 deleterious, as well as presumably deleterious variants, some of which reside in genes
274 with known disease associations, again suggesting moderate to strong purifying
275 selection. The dependency of allele frequencies on the proximity to classical HLA genes
276 supports the notion that balancing selection acts on neighboring variation via linkage
277 disequilibrium. Independent support for our findings comes also from two earlier studies
278 that described elevated levels of genetic diversity around specific MHC genes and
279 suggested that this might be maintained by balancing selection on the given neighboring
280 MHC locus (Horton et al. 1998; Shiina et al. 2006).

281
282 The observed reduction in the number of polymorphic sites in regions close to loci under
283 multi-locus balancing selection both in simulations and empirical data is an interesting
284 observation, given the expectation that balancing selection generally increases genetic
285 diversity around the loci under selection (Navarro and Barton 2002; Charlesworth 2006).
286 This is also a significant difference to the effect of genetic hitchhiking after selective
287 sweeps, which has been shown to reduce the number of neutral variants in neighboring
288 regions but not to affect the number of deleterious variants (Chun and Fay 2011).
289 However, the unexpected reduction in the number of polymorphic sites with balancing
290 selection at the neighboring HLA is likely due to the recurrent occurrence of new
291 mutations at one of the balanced HLA sites, a central aspect of our multi-locus
292 simulation model. Such a new variant, being subject to overdominant selection, would
293 drive the haplotype on which it occurred to higher frequency (but not to fixation),
294 dragging along any linked neighboring variant. At the same time the rising frequency of
295 this haplotype would replace other haplotypes that don't carry the new balanced variant.
296 Consequently, a polymorphic site linked to a balanced locus can have two fates: Either

297 one of its minor alleles is present on the new balanced haplotype, so that its frequency
298 rises together with the balanced haplotype, or the new balanced haplotype does not
299 carry the site's minor alleles, so that its frequency declines. For very low frequency
300 variants (e.g. singletons) the latter fate will potentially lead to a complete loss of the
301 minor allele at this neighboring site, and thus to a complete loss of the segregating site.
302 Together these two mechanisms lead to the signature that we observe: Less
303 polymorphic sites around the HLA gene but a higher derived allele frequency at those
304 sites that remain in the population.

305 Such a scenario of recurrent but incomplete sweeps could indeed result from several
306 mechanisms of balancing selection, such as negative frequency-dependent selection
307 and fluctuating selection in time and space, which are also thought to affect the HLA
308 (Meyer and Thomson 2001; Spurgin and Richardson 2010). Previous theoretical work
309 has even suggested that symmetrical overdominance alone is not sufficient to explain
310 the entire allelic diversity seen at the HLA (De Boer et al. 2004). Furthermore, the HLA
311 exhibits identity-by-descent patterns that cannot be explained by ancient balancing
312 selection, suggesting ongoing selection and very recent selective events (Albrechtsen et
313 al. 2010), for instance through local adaptation after human migration and/or drastic
314 epidemic events in recent human history (Zhou et al. 2016). All these processes are
315 likely to contribute to the distribution of linked deleterious variation around the classical
316 HLA genes.

317
318 Mechanistically, it had also been proposed that the excessive heterozygosity in the MHC
319 region, caused by balancing selection, may reduce the phenotypic expression of
320 recessive deleterious variants, rendering purging mechanisms less effective and thus
321 allowing for the accumulation of a recessive 'sheltered load' that could ultimately even
322 contribute to inbreeding depression (Charlesworth and Willis 2009). Such a scenario
323 was first invoked to explain the unexpectedly long terminal branches in genealogies of
324 genes under balancing selection (Uyenoyama 1997; for a schematic see van Oosterhout
325 2009) and has since been investigated by theoretical and empirical studies on the self-
326 incompatibility determining S locus, the most compelling example for balancing selection
327 in plants. Those studies suggested that recessive deleterious mutations may accumulate

328 around the S locus, presumably protected from purifying selection by the excessive
329 heterozygosity in this region, and, once established, contribute to the disadvantage of
330 homozygotes (Stone 2004; Llaurens et al. 2009). Such a 'sheltered load' mechanism
331 has later also been proposed to contribute to the diversity found in the MHC region in
332 vertebrates (van Oosterhout 2009; Llaurens et al. 2012). However, in contrast to
333 previous work, our models are based on empirical parameter values, showing an
334 accumulating deleterious load under realistic levels of recombination. In addition we
335 here provide direct empirical evidence of elevated deleterious variation in the human
336 MHC region. While it is conceivable that some of these observed deleterious variants
337 represent such a 'sheltered load', our simulations assume additive (and not recessive)
338 effects for all non-HLA variants, indicating that there must be other mechanisms besides
339 the 'heterozygosity shelter' for recessive variants that contribute to the observed
340 deleterious load.

341 Contributing to the observed excess of deleterious variants could also be a local
342 alteration of the effective population size, and consequently of the efficacy of selection.
343 Balancing selection, and specifically symmetrical overdominance, can lead to a limited
344 number of dominating haplotypes that are stably balanced at intermediate frequencies in
345 the population. Without recombination, such dominating haplotypes could be seen as
346 essentially subdividing the population in this genomic region (Charlesworth 2006), so
347 that genetic drift of linked variants occurs only within the same haplotype background.
348 This would result in a locally reduced effective population size per haplotype, rendering
349 selection less efficient within a given haplotype background, even though, when
350 estimated across all haplotypes, balancing selection increases the effective size of the
351 total population (Charlesworth 2009). This phenomenon is strongly based on genetic
352 linkage, and with recombination the effect on linked variants will quickly break down with
353 increasing distance to the balanced locus (Hudson and Kaplan 1988). Such a
354 dependence between the excess deleterious load and the rate of recombination could
355 indeed be observed both in the simulations and in the empirical data (here using
356 physical distance along the chromosome as a proxy for genetic linkage).

357

358 Overall, our results show that balancing selection has the potential to maintain
359 deleterious genetic variants at considerable frequencies in natural populations,
360 especially if they are only moderately deleterious, e.g. causing late-onset autoimmunity.
361 The extent of the deleterious mutational load on a given MHC haplotype will ultimately
362 depend on its aggregate detrimental effect and the role of linked HLA alleles, resulting
363 from an evolutionary trade-off between the antagonistic effects of the number (and
364 severity) of accumulating deleterious variants and the selective advantage of individual
365 allelic variability at the classical HLA genes. Interestingly, more divergent HLA
366 genotypes are assumed to confer broader immune-surveillance against antigens (Lenz
367 2011), and a positive correlation between the pathogen diversity in a given population
368 and its HLA allele pool diversity has been reported (Prugnolle et al. 2005). It could thus
369 be hypothesized that environments with higher pathogen diversity, and thus stronger
370 selection for high individual allelic diversity at HLA, may allow for the maintenance of
371 MHC haplotypes with a larger deleterious load (Dean et al. 2002). Once the selective
372 pathogenic pressure declines (either through host migration to a less pathogenic
373 environment or through environmental changes), this could then result in increased
374 expression of phenotypes associated with the deleterious load. Indeed, first evidence
375 associating heterozygosity at the HLA with increased risk for some autoimmune
376 diseases has been reported (Lenz et al. 2015), but further analyses are required to
377 evaluate this hypothesis.

378
379 The observed effect of balancing selection on linked variation is evidently dependent on
380 recombination rate around the locus under balancing selection (Hudson and Kaplan
381 1988; Barton 2010). The MHC region is in fact known to harbor a substantial number of
382 recombination hot-spots (de Bakker et al. 2006), which might have evolved due to the
383 above described dynamics. On the other hand, the MHC region also exhibits even more
384 extreme cases of multi-locus balancing selection than modeled here, where LD spans
385 across separate HLA loci. For instance, the well-documented COX haplotype covers the
386 entire MHC region and has a population frequency of about 10% in Northern Europe
387 (Stewart et al. 2004). Epistasis among different HLA genes, for instance due to
388 advantageous allele combinations, may select against recombination and maintain such

389 long-range haplotypes (Penman et al. 2013). The genome-wide extent of the observed
390 effect remains to be investigated, but given the increasing number of regions detected to
391 evolve under balancing selection (Leffler et al. 2013), there is reason to expect that this
392 mechanism may significantly contribute to the prevalence of heritable human diseases.

393

394

395 **Materials and Methods**

396 **GWAS summary**

397 Data for GWAS hit summary was downloaded from the NHGRI GWAS catalog (available
398 at: <http://www.genome.gov/gwastudies>, accessed [02/15/2015]) (Welter et al. 2014).

399 This catalog represents a comprehensive collection of published genome-wide
400 association studies (currently 1,755 studies), including the investigated trait and the
401 chromosomal location of one or more (median: 4) independent associations per study.

402 For some traits the data contains multiple studies (top five: type 2 diabetes: 37, breast
403 cancer: 26, schizophrenia: 22, body mass index: 21, height: 21; median number of
404 studies per trait: 1), but preferentially lists novel associations to avoid redundancy. While
405 this data may suffer from a bias towards over-studied traits, it also reflects the frequency
406 and thus importance of a given trait in the population. We focused only on autosomal
407 associations and chromosome length was standardized to 30 location bins per
408 chromosome for better visualization.

409

410 **Simulating multi-locus balancing selection**

411 In order to test whether balancing selection on a particular locus can lead to a shift in the
412 site frequency spectrum of adjacent genomic regions and prevent deleterious mutations
413 from being purged by purifying selection, we employed a forward simulation approach.

414 To this end, we developed the program Forward Simulation (available at
415 <http://forwardsimulation.sourceforge.net>), which is based on the Wright-Fisher model
416 and allows to define distinct selection regimes for separate regions or sites of the
417 simulated genome. The program uses a multiplicative fitness framework with diallelic
418 sites, with fitness $f = (1+sh)^n$, where s and h are the selection and dominance
419 coefficients per site, respectively, and n is the number of sites in the genome. A negative

420 selection coefficient s leads to negative (purifying) selection and vice versa. The
421 simulated genome contained a virtual HLA gene of the median genomic length of
422 classical HLA genes (start of first exon to end of last exon, including introns; 5,385 bp).
423 This HLA gene was flanked on both sides by surrounding regions (set to evolve under
424 different levels of purifying selection). The length of these two regions equaled half the
425 median distance between adjacent classical HLA genes in the MHC region (total length:
426 37,122 bp). Following the multi-locus balancing selection scenario by Navarro and
427 Barton (2002), a number of polymorphic sites along the virtual HLA gene were specified
428 to evolve under the HLA selection regime at the beginning of each simulation (after
429 burnin). Once the sites had been specified, they were free to maintain, lose, or regain
430 polymorphism throughout the simulation. These sites mimic the functional polymorphism
431 seen in classical HLA genes, where sites along the exons coding for the antigen binding
432 groove of the HLA molecule show variation that is thought to evolve under balancing
433 selection (Reche and Reinherz 2003; Furlong and Yang 2008). See **supplementary fig.**
434 **S5** (Supplementary Material online) for a schematic of the simulated genome. For
435 exploring the effect of balancing selection on surrounding regions in a near-realistic
436 scenario of MHC evolution, we used empirical values from the literature for all fixed
437 parameters: effective population size $N_e = 10,000$ (Takahata et al. 1995), mutation rate μ
438 $= 1.38e-8$ (Scally and Durbin 2012), recombination rate $r = 0.44$ cM/Mb (de Bakker et al.
439 2006; Taylan and Altıok 2012). Average selection at the classical HLA genes has been
440 estimated as $s = 0.013$ (Satta et al. 1994; Slatkin and Muirhead 2000; Yasukochi and
441 Satta 2013). Following Navarro and Barton (2002), we simulate multi-locus balancing
442 selection as symmetrical overdominance, setting the selection coefficient s_{HLA} for all
443 sites with HLA selection regime to 0.013 for heterozygotes and to 0 for homozygotes.
444 For simulation scenarios with no selection at HLA, s_{HLA} was set to 0. The scenarios with
445 recurrent selective sweeps at the HLA gene were simulated by setting $s_{HLA} = 0.013$, with
446 a dominance coefficient $h = 0.5$, and reverting sites with fixed derived allele back to the
447 ancestral allele in order to allow for sweeps to reoccur. The number of sites to evolve
448 under the HLA selection regime was 117, which corresponds to the average number of
449 amino acid residues coding for the peptide binding domain of the classical HLA
450 molecules. All other sites along the entire genome were set to evolve either neutrally (s

451 at non-HLA or 'neighboring' sites: $s_{nb} = 0$) or under purifying selection, with s_{nb} ranging
452 from -0.0001 to -0.01, and dominance $h = 0.5$ (additive). The simulations were run for
453 150,000 generations, after an initial burnin of 100,000 generations intended to generate
454 an equilibrium level of neutral variation ($\pi \approx 4N\mu$; **supplementary fig. S6**,
455 Supplementary Material online). During the burnin period, selection at the HLA gene was
456 turned off ($s_{HLA} = 0$), so that all sites evolved according to the specified level of purifying
457 selection ($s_{nb} = 0$ to -0.0001), including the possibility for fixation of derived alleles. Each
458 simulated scenario was replicated 100 times. Reported results are based only on
459 variants in the surrounding region, thus excluding the virtual HLA gene.

460

461 **Exome sequencing data and MHC region**

462 Genome-wide coding sequence variation data of about 6,500 human individuals was
463 obtained from the Exome Variant Server (NHLBI GO Exome Sequencing Project
464 (ESP6500 release), Seattle, WA; <http://evs.gs.washington.edu/EVS/>; accessed 11/2012)
465 and parsed with custom scripts. Only variants that passed the internal EVS quality filter
466 were recorded and their number and frequencies averaged per gene to account for
467 differences in cds length among genes. Missense variants were further split into
468 functional categories, as determined by Polyphen-2 (Adzhubey et al. 2010). All
469 functional annotations used here are based on the non-reference allele. For simplicity
470 we subsequently call this the derived allele, even though the true ancestral state may
471 not in all cases be reliably resolved. However, our analyses rely only on the functional
472 annotation of a given allele and are independent of the true ancestral/derived state.

473

474 The 3.5 Mb region of the classical major histocompatibility complex (MHC) on
475 chromosome 6 was defined following Horton *et al.* (2004), starting with the gene *ZFP57*
476 at position chr6:29640169 and ending with the gene *HCG24* at position chr6:33115544.
477 This excludes the extended MHC regions which contain large numbers of olfactory
478 receptor and histone genes that are likely subject to different types of selection regimes
479 and are thus outside the scope of this study. Within the classical MHC region, the
480 ESP6500 data provided sequence variation data for 124 genes, which corresponds to
481 73.4% of the known protein-coding genes in this region. This fraction corresponds well

482 with the genome-wide proportion (73.5%) of genes sufficiently covered by the ESP data.
483 Eleven of the 124 genes covered in the MHC region are classical HLA genes (class I:
484 *HLA-A, -B, -C*; class II: *HLA-DPA1, -DPB1, -DQA1, -DQA2, -DQB1, -DRA, -DRB1, -*
485 *DRB5*).

486
487 Site frequency spectra (SFS), based on derived allele frequencies, were obtained for
488 each of the variant categories across the entire exome as well as for only the MHC
489 region. SFS were compared using the non-parametric Wilcoxon rank-sum test. Within
490 the MHC region, we also calculated for each variant the distance in base pairs from the
491 nearest classical HLA locus and correlated this distance with the derived allele
492 frequency using Spearman correlation. We furthermore compared observed parameter
493 values in the MHC region with genome-wide expectations by Monte Carlo sampling
494 (1,000 replications) from all covered genes in the ESP6500 data. Data analysis was
495 done in R (ver. 3.1.2) (R Development Core Team 2014).

496
497 Disease associations for genes in the MHC region with the highest average frequency of
498 potentially deleterious variants were obtained from NHGRI GWAS catalog (see above)
499 and the NIH Genetic Association Database (available at:
500 <http://geneticassociationdb.nih.gov>, accessed [12/14/2013]) (Becker et al. 2004).

501
502

503 **Acknowledgements:**

504 We thank D. Balick for fruitful discussions and comments on a previous version of the
505 manuscript. P. Polak kindly shared gene-wise mutation rate data. We are also grateful to
506 the constructive and thoughtful comments of the editor and three reviewers that helped
507 to improve the manuscript. This work was supported by German Research Foundation
508 (DFG) grants LE 2593/1-1 and LE 2593/2-1 (to T.L.L.), and National Institutes of Health
509 (NIH) grants R01 GM078598 and R01 MH101244 (to S.R.S. and D.M.J.).

510
511
512

513 **References:**

- 514 Adzhubey IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov
515 AS, Sunyaev SR. 2010. A method and server for predicting damaging missense
516 mutations. *Nat Methods* 7:248-249.
- 517 Albrechtsen A, Moltke I, Nielsen R. 2010. Natural Selection and the Distribution of
518 Identity-by-Descent in the Human Genome. *Genetics* 186:295-308.
- 519 Andrés AM, Hubisz MJ, Indap A, Torgerson DG, Degenhardt JD, Boyko AR, Gutenkunst
520 RN, White TJ, Green ED, Bustamante CD, et al. 2009. Targets of balancing
521 selection in the human genome. *Mol Biol Evol* 26:2755-2764.
- 522 Assaf ZJ, Petrov DA, Blundell JR. 2015. Obstruction of adaptation in diploids by
523 recessive, strongly deleterious alleles. *Proc Natl Acad Sci USA* 112:E2658-
524 E2666.
- 525 Asthana S, Schmidt S, Sunyaev SR. 2005. A limited role for balancing selection. *Trends*
526 *Genet* 21:30-32.
- 527 Barton NH. 2010. Genetic linkage and natural selection. *Philos Trans R Soc B Biol Sci*
528 365:2559-2569.
- 529 Becker KG, Barnes KC, Bright TJ, Wang SA. 2004. The Genetic Association Database.
530 *Nat Genet* 36:431-432.
- 531 Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE,
532 Adams MD, Schmidt S, Sninsky JJ, Sunyaev SR, et al. 2008. Assessing the
533 Evolutionary Impact of Amino Acid Mutations in the Human Genome. *PLoS Genet*
534 4:e1000083.
- 535 Bubb KL, Bovee D, Buckley D, Haugen E, Kibukawa M, Paddock M, Palmieri A,
536 Subramanian S, Zhou Y, Kaul R, et al. 2006. Scan of human genome reveals no
537 new loci under ancient balancing selection. *Genetics* 173:2165-2177.
- 538 Bustamante CD, Fledel-Alon A, Williamson S, Nielsen R, Todd Hubisz M, Glanowski S,
539 Tanenbaum DM, White TJ, Sninsky JJ, Hernandez RD, et al. 2005. Natural
540 selection on protein-coding genes in the human genome. *Nature* 437:1153-1157.
- 541 Charlesworth B. 2009. Fundamental concepts in genetics: Effective population size and
542 patterns of molecular evolution and variation. *Nat Rev Genet* 10:195-205.

- 543 Charlesworth B, Charlesworth D, Barton NH. 2003. The effects of genetic and
544 geographic structure on neutral variation. *Annu Rev Ecol Evol Syst* 34:99-125.
- 545 Charlesworth B, Morgan MT, Charlesworth D. 1993. The effect of deleterious mutations
546 on neutral molecular variation. *Genetics* 134:1289-1303.
- 547 Charlesworth D. 2006. Balancing selection and its effects on sequences in nearby
548 genome regions. *PLoS Genet* 2:e64.
- 549 Charlesworth D, Willis JH. 2009. The genetics of inbreeding depression. *Nat Rev Genet*
550 10:783-796.
- 551 Chun S, Fay JC. 2011. Evidence for hitchhiking of deleterious mutations within the
552 human genome. *PLoS Genet* 7:e1002240.
- 553 Connallon T, Clark AG. 2013. Antagonistic versus nonantagonistic models of balancing
554 selection: Characterizing the relative timescales and hitchhiking effects of partial
555 selective sweeps. *Evolution* 67:908-917.
- 556 de Bakker PIW, McVean G, Sabeti PC, Miretti MM, Green T, Marchini J, Ke X, Monsuur
557 AJ, Whittaker P, Delgado M, et al. 2006. A high-resolution HLA and SNP
558 haplotype map for disease association studies in the extended human MHC. *Nat*
559 *Genet* 38:1166-1172.
- 560 De Boer RJ, Borghans JA, van Boven M, Kesmir C, Weissing FJ. 2004. Heterozygote
561 advantage fails to explain the high degree of polymorphism of the MHC.
562 *Immunogenetics* 55:725-731.
- 563 Dean M, Carrington M, O'Brien SJ. 2002. Balanced polymorphism selected by genetic
564 versus infectious human disease. *Annu Rev Genomics Hum Genet* 3:263-292.
- 565 DeGiorgio M, Lohmueller KE, Nielsen R. 2014. A Model-Based Approach for Identifying
566 Signatures of Ancient Balancing Selection in Genetic Data. *PLoS Genet*
567 10:e1004561.
- 568 Dudley JT, Kim Y, Liu L, Markov GJ, Gerold K, Chen R, Butte AJ, Kumar S. 2012.
569 Human genomic disease variants: A neutral evolutionary explanation. *Genome*
570 *Res* 22:1383-1394.
- 571 Fay JC, Wu C-I. 2000. Hitchhiking under positive darwinian selection. *Genetics*
572 155:1405-1413.

- 573 Francioli LC, Polak PP, Koren A, Menelaou A, Chun S, Renkens I, Genome of the
574 Netherlands Consortium, Duijn CMv, Swertz M, Wijmenga C, et al. 2015.
575 Genome-wide patterns and properties of de novo mutations in humans. *Nat*
576 *Genet* 47:822-826.
- 577 Fu W, O'Connor TD, Jun G, Kang HM, Abecasis G, Leal SM, Gabriel S, Rieder MJ,
578 Altshuler D, Shendure J, et al. 2013. Analysis of 6,515 exomes reveals the recent
579 origin of most human protein-coding variants. *Nature* 493:216-220.
- 580 Fumagalli M, Cagliani R, Pozzoli U, Riva S, Comi GP, Menozzi G, Bresolin N, Sironi M.
581 2009. Widespread balancing selection and pathogen-driven selection at blood
582 group antigen genes. *Genome Res* 19:199-212.
- 583 Furlong R, Yang Z. 2008. Diversifying and purifying selection in the peptide binding
584 region of DRB in mammals. *J Mol Evol* 66:384-394.
- 585 Gao Z, Przeworski M, Sella G. 2015. Footprints of ancient balanced polymorphisms in
586 genetic variation data from closely related species. *Evolution* 69:431-446.
- 587 Gillespie JH. 2000. Genetic drift in an infinite population: The pseudohitchhiking model.
588 *Genetics* 155:909-919.
- 589 Gokcumen O, Zhu Q, Mulder LCF, Iskow RC, Austermann C, Scharer CD, Raj T, Boss
590 JM, Sunyaev SR, Price A, et al. 2013. Balancing selection on a regulatory region
591 exhibiting ancient variation that predates human–Neandertal divergence. *PLoS*
592 *Genet* 9:e1003404.
- 593 Grimsley C, Mather KA, Ober C. 1998. HLA-H: a pseudogene with increased variation
594 due to balancing selection at neighboring loci. *Mol Biol Evol* 15:1581-1588.
- 595 Hartfield M, Otto SP. 2011. Recombination and hitchhiking of deleterious alleles.
596 *Evolution* 65:2421-2434.
- 597 Henn BM, Botigue LR, Bustamante CD, Clark AG, Gravel S. 2015. Estimating the
598 mutation load in human genomes. *Nat Rev Genet* 16:333-343.
- 599 Horton R, Niblett D, Milne S, Palmer S, Tubby B, Trowsdale J, Beck S. 1998. Large-
600 scale sequence comparisons reveal unusually high levels of variation in the HLA-
601 DQB1 locus in the class II region of the human MHC. *J Mol Biol* 282:71-97.

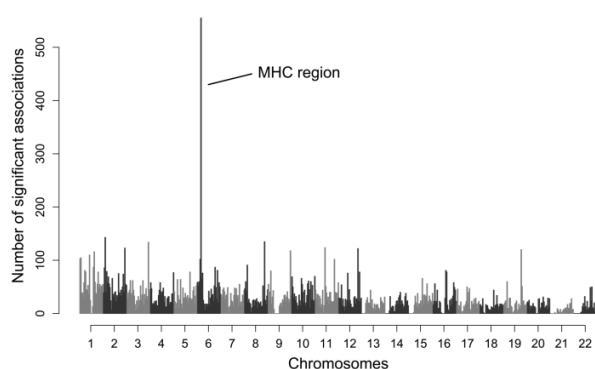
- 602 Horton R, Wilming L, Rand V, Lovering RC, Bruford EA, Khodiyar VK, Lush MJ, Povey
603 S, Talbot CC, Wright MW, et al. 2004. Gene map of the extended human MHC.
604 Nat Rev Genet 5:889-899.
- 605 Hudson RR, Kaplan NL. 1988. The coalescent process in models with selection and
606 recombination. Genetics 120:831-840.
- 607 Kaplan NL, Darden T, Hudson RR. 1988. The coalescent process in models with
608 selection. Genetics 120:819-829.
- 609 Kiezun A, Pulit SL, Francioli LC, van Dijk F, Swertz M, Boomsma DI, van Duijn CM,
610 Slagboom PE, van Ommen GJB, Wijmenga C, et al. 2013. Deleterious Alleles in
611 the Human Genome Are on Average Younger Than Neutral Alleles of the Same
612 Frequency. PLoS Genet 9:e1003301.
- 613 Kryukov GV, Pennacchio LA, Sunyaev SR. 2007. Most Rare Missense Alleles Are
614 Deleterious in Humans: Implications for Complex Disease and Association
615 Studies. Am J Hum Genet 80:727-739.
- 616 Kryukov GV, Shpunt A, Stamatoyannopoulos JA, Sunyaev SR. 2009. Power of deep, all-
617 exon resequencing for discovery of human trait genes. Proc Natl Acad Sci USA
618 106:3871-3876.
- 619 Leffler EM, Gao Z, Pfeifer S, Ségurel L, Auton A, Venn O, Bowden R, Bontrop R, Wall
620 JD, Sella G, et al. 2013. Multiple instances of ancient balancing selection shared
621 between humans and chimpanzees. Science 339:1578-1582.
- 622 Lenz TL. 2011. Computational prediction of MHC II-antigen binding supports divergent
623 allele advantage and explains trans-species polymorphism. Evolution 65:2380-
624 2390.
- 625 Lenz TL, Deutsch AJ, Han B, Hu X, Okada Y, Eyre S, Zhernakova A, Huizinga TWJ,
626 Abecasis G, Becker J, et al. 2015. Widespread non-additive and interaction
627 effects within HLA loci modulate the risk of autoimmune diseases. Nat Genet
628 47:1085-1090.
- 629 Lewontin RC, Kojima K-i. 1960. The Evolutionary Dynamics of Complex Polymorphisms.
630 Evolution 14:458-472.

- 631 Llaurens V, Gonthier L, Billiard S. 2009. The sheltered genetic load linked to the S locus
632 in plants: new insights from theoretical and empirical approaches in sporophytic
633 self-incompatibility. *Genetics* 183:1105-1118.
- 634 Llaurens V, McMullan M, van Oosterhout C. 2012. Cryptic MHC polymorphism revealed
635 but not explained by selection on the class IIB peptide-binding region. *Mol Biol*
636 *Evol* 29:1631-1644.
- 637 Marsden CD, Ortega-Del Vecchyo D, O'Brien DP, Taylor JF, Ramirez O, Vilà C,
638 Marques-Bonet T, Schnabel RD, Wayne RK, Lohmueller KE. 2016. Bottlenecks
639 and selective sweeps during domestication have increased deleterious genetic
640 variation in dogs. *Proc Natl Acad Sci USA* 113:152-157.
- 641 Maurano MT, Humbert R, Rynes E, Thurman RE, Haugen E, Wang H, Reynolds AP,
642 Sandstrom R, Qu H, Brody J, et al. 2012. Systematic Localization of Common
643 Disease-Associated Variation in Regulatory DNA. *Science* 337:1190-1195.
- 644 Maynard Smith J, Haigh J. 1974. The hitch-hiking effect of a favourable gene. *Genet*
645 *Res* 23:23-35.
- 646 Meyer D, Thomson G. 2001. How selection shapes variation of the human major
647 histocompatibility complex: a review. *Ann Hum Genet* 65:1-26.
- 648 Navarro A, Barton NH. 2002. The effects of multilocus balancing selection on neutral
649 variability. *Genetics* 161:849-863.
- 650 Neefjes J, Jongsma MLM, Paul P, Bakke O. 2011. Towards a systems understanding of
651 MHC class I and MHC class II antigen presentation. *Nat Rev Immunol* 11:823-
652 836.
- 653 O'hUigin C, Satta Y, Hausmann A, Dawkins RL, Klein J. 2000. The implications of
654 intergenic polymorphism for major histocompatibility complex evolution. *Genetics*
655 156:867-877.
- 656 Penman BS, Ashby B, Buckee CO, Gupta S. 2013. Pathogen selection drives
657 nonoverlapping associations between HLA loci. *Proc Natl Acad Sci USA*
658 110:19645-19650.
- 659 Prugnolle F, Manica A, Charpentier M, Guegan JF, Guernier V, Balloux F. 2005.
660 Pathogen-driven selection and worldwide HLA class I diversity. *Curr Biol* 15:1022-
661 1027.

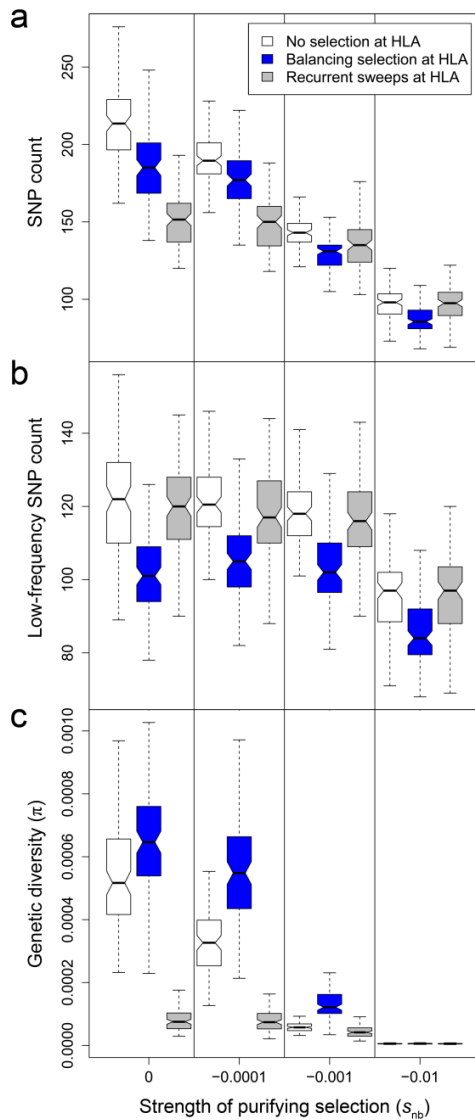
- 662 R Development Core Team. 2014. R: A language and environment for statistical
663 computing. Version Version 3.1.2. Vienna, Austria: R Foundation for Statistical
664 Computing.
- 665 Reche PA, Reinherz EL. 2003. Sequence variability analysis of human class I and class
666 II MHC molecules: Functional and structural correlates of amino acid
667 polymorphisms. *J Mol Biol* 331:623-641.
- 668 Satta Y, O'hUigin C, Takahata N, Klein J. 1994. Intensity of natural selection at the
669 major histocompatibility complex loci. *Proc Natl Acad Sci USA* 91:7184-7188.
- 670 Scally A, Durbin R. 2012. Revising the human mutation rate: implications for
671 understanding human evolution. *Nat Rev Genet* 13:745-753.
- 672 Sellis D, Callahan BJ, Petrov DA, Messer PW. 2011. Heterozygote advantage as a
673 natural consequence of adaptation in diploids. *Proc Natl Acad Sci USA*
674 108:20666-20671.
- 675 Shiina T, Ota M, Shimizu S, Katsuyama Y, Hashimoto N, Takasu M, Anzai T, Kulski JK,
676 Kikkawa E, Naruse T, et al. 2006. Rapid evolution of major histocompatibility
677 complex class I genes in primates generates new disease alleles in humans via
678 hitchhiking diversity. *Genetics* 173:1555-1570.
- 679 Slatkin M, Muirhead CA. 2000. A method for estimating the intensity of overdominant
680 selection from the distribution of allele frequencies. *Genetics* 156:2119-2126.
- 681 Spurgin LG, Richardson DS. 2010. How pathogens drive genetic diversity: MHC,
682 mechanisms and misunderstandings. *Proc R Soc B Biol Sci* 277:979-988.
- 683 Stewart CA, Horton R, Allcock RJN, Ashurst JL, Atrazhev AM, Coggill P, Dunham I,
684 Forbes S, Halls K, Howson JMM, et al. 2004. Complete MHC haplotype
685 sequencing for common disease gene mapping. *Genome Res* 14:1176-1187.
- 686 Stone JL. 2004. Sheltered load associated with S-alleles in *Solanum carolinense*.
687 *Heredity* 92:335-342.
- 688 Takahata N, Satta Y, Klein J. 1995. Divergence time and population size in the lineage
689 leading to modern humans. *Theor Popul Biol* 48:198-221.
- 690 Taylan F, Altioek E. 2012. Meiotic recombinations within major histocompatibility complex
691 of human embryos. *Immunogenetics* 64:839-844.

- 692 Teixeira JC, de Filippo C, Weihmann A, Meneu JR, Racimo F, Dannemann M, Nickel B,
693 Fischer A, Halbwax M, Andre C, et al. 2015. Long-term balancing selection in
694 LAD1 maintains a missense trans-species polymorphism in humans,
695 chimpanzees, and bonobos. *Mol Biol Evol* 32:1186-1196.
- 696 Trowsdale J. 2011. The MHC, disease and selection. *Immunol Lett* 137:1-8.
- 697 Uyenoyama MK. 1997. Genealogical structure among alleles regulating self-
698 incompatibility in natural populations of flowering plants. *Genetics* 147:1389-1400.
- 699 van Oosterhout C. 2009. A new theory of MHC evolution: beyond selection on the
700 immune genes. *Proc R Soc B Biol Sci* 276:657-665.
- 701 Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P,
702 Manolio T, Hindorff L, et al. 2014. The NHGRI GWAS Catalog, a curated
703 resource of SNP-trait associations. *Nucleic Acids Res* 42:D1001-D1006.
- 704 Yasukochi Y, Satta Y. 2013. Current perspectives on the intensity of natural selection of
705 MHC loci. *Immunogenetics* 65:479-483.
- 706 Zhou Q, Zhao L, Guan Y. 2016. Strong Selection at MHC in Mexicans since Admixture.
707 *PLoS Genet* 12:e1005847.
- 708

709 **Figures**



710
711 **Figure 1: Number of significant GWAS associations along the genome.** The
712 chromosomal location of significant trait associations from genome-wide association
713 studies (GWAS, N=18,682) are shown for all autosomes. Data from NHGRI GWAS
714 catalog.
715



716

717 **Figure 2. Simulated polymorphism in regions surrounding the HLA gene under**

718 **different selection scenarios.** Polymorphic sites (SNPs) along the regions around an

719 HLA gene are derived from simulations with three different selection scenarios on the

720 HLA gene (white: no selection on HLA, blue: balancing selection, grey: recurrent sweeps

721 of positive selection). Standard box plots show the median number of (a) all SNPs, (b)

722 only SNPs with derived allele frequency < 0.01, and (c) the genetic diversity (π) across

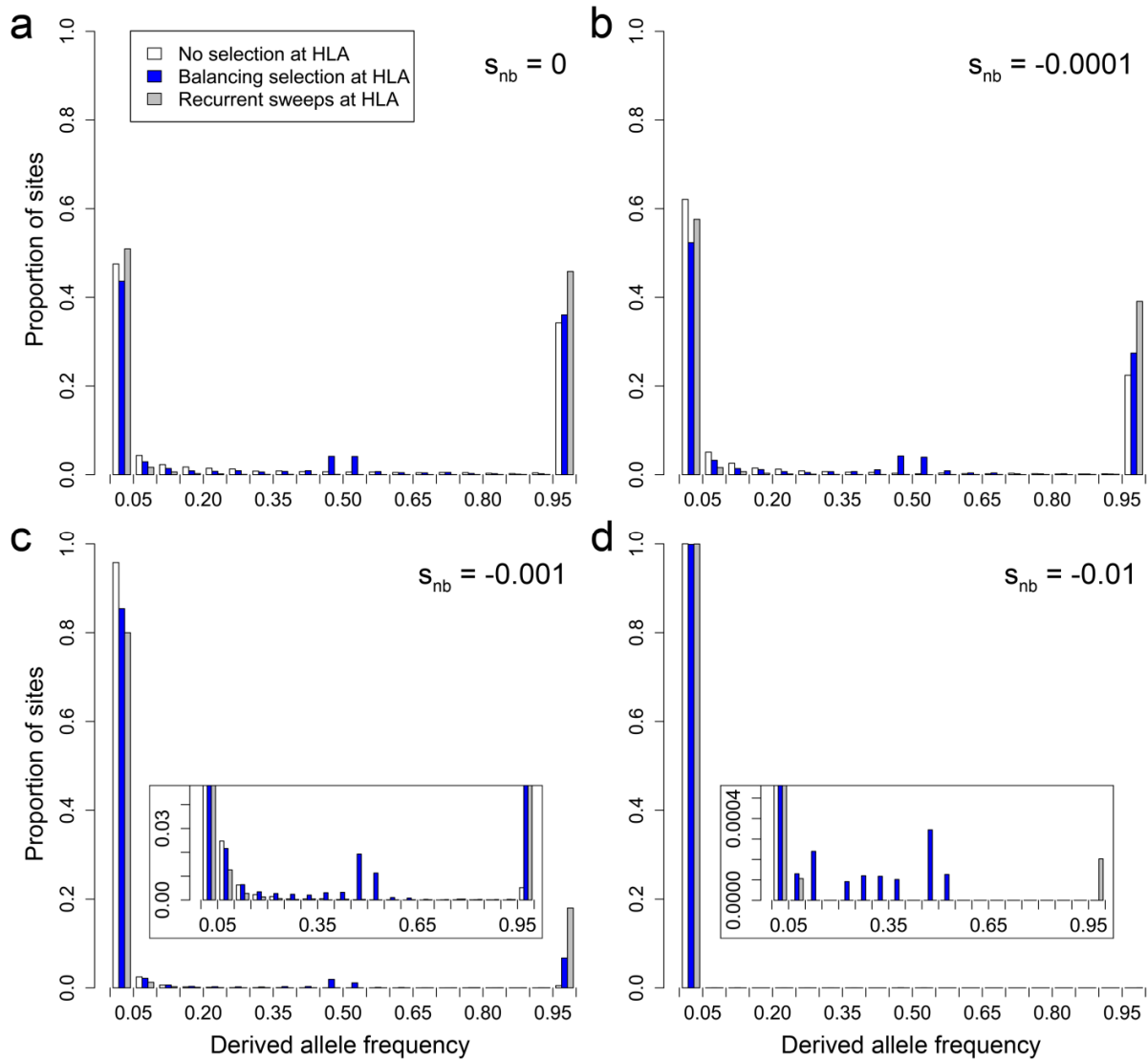
723 all sites. Variants in surrounding regions evolved neutrally ($s_{nb} = 0$) or under co-dominant

724 purifying selection with $s_{nb} = -0.0001$, $s_{nb} = -0.001$, or $s_{nb} = -0.01$, respectively. Non-

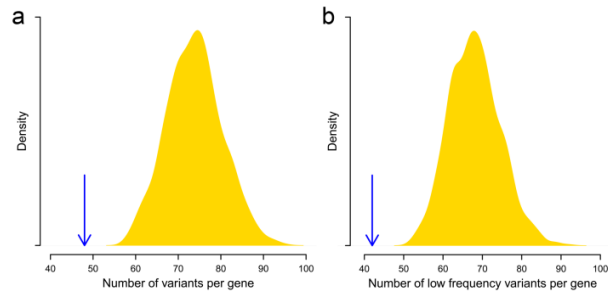
725 overlapping notches between box plots indicate significant difference. Note the different

726 y-axis scales.

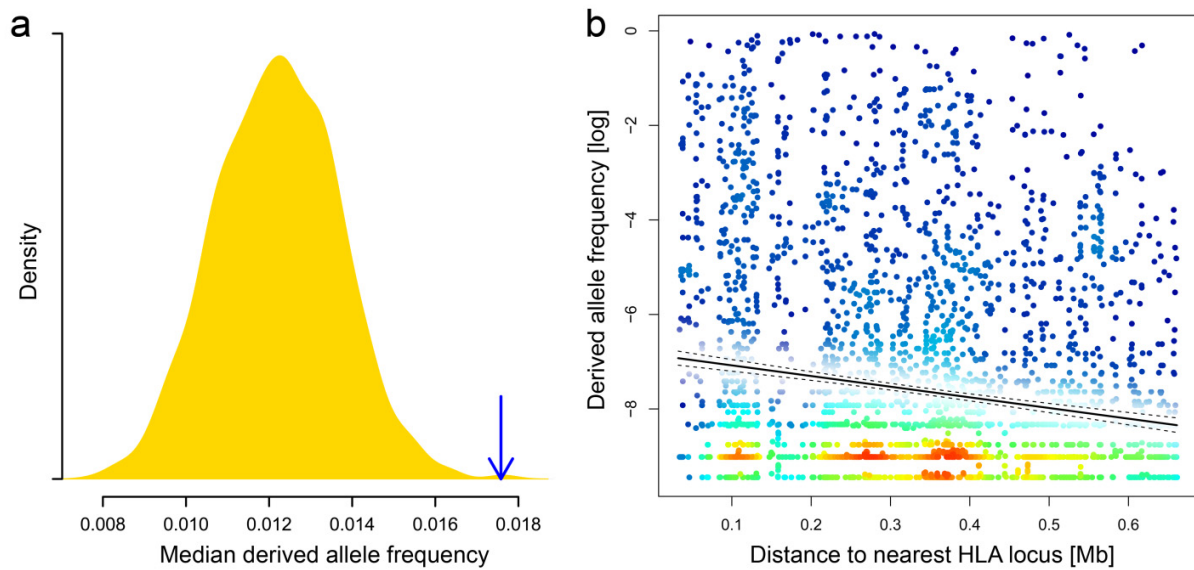
727



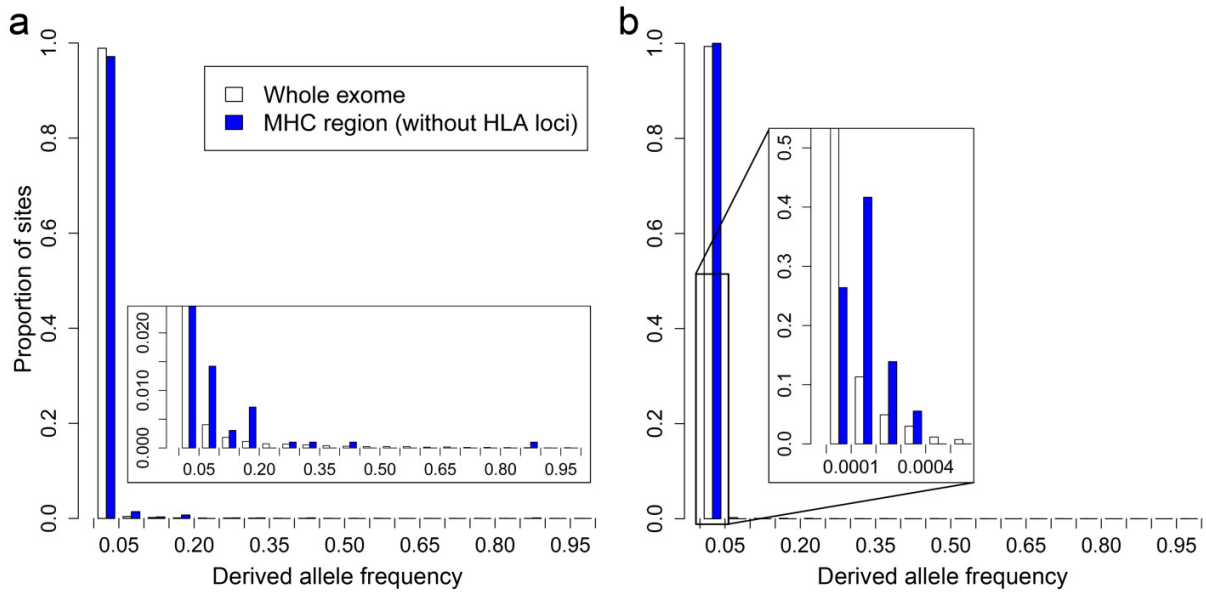
728
729 **Figure 3. Simulated site frequency spectrum of variants surrounding the HLA**
730 **gene under different selection scenarios.** Derived allele frequencies along the
731 regions around an HLA gene are derived from simulations with three different selection
732 scenarios on the HLA gene (white: no selection on HLA, blue: balancing selection, grey:
733 recurrent sweeps of positive selection). Variants in neighboring regions evolved (a)
734 neutrally ($s_{nb} = 0$) or under co-dominant purifying selection with (b) $s_{nb} = -0.0001$,
735 (c) $s_{nb} = -0.001$, or (d) $s_{nb} = -0.01$, respectively. Note the different y-axis scales in the
736 zoomed insets of panels (c) and (d) for better visualization.
737



738
739 **Figure 4. Observed average number of polymorphic sites per gene.** The median
740 number of polymorphic sites per gene is averaged over the 113 genes represented in
741 the MHC region (blue arrow, excluding classical HLA genes). Also shown is the
742 distribution of equivalent values for 1,000 Monte-Carlo-sampled sets of 113 random
743 genes from the entire exome. Represented are (a) all SNPs and (b) only SNPs with
744 derived allele frequency < 0.01.
745



746
747 **Figure 5. Observed distribution of derived allele frequencies.** Derived alleles are
748 defined as the non-reference allele at polymorphic sites in the Exome Sequencing
749 Project data. **(a)** The median frequency of derived alleles per gene is averaged over the
750 113 genes represented in the MHC region (blue arrow, excluding classical HLA genes).
751 The yellow density curve shows the distribution of equivalent values for 1000 Monte-
752 Carlo-sampled sets of 113 random genes from the entire exome. **(b)** Derived allele
753 frequency at polymorphic sites (N = 4,175) in the MHC region is shown in relation to
754 distance to the nearest classical HLA locus. Solid and dashed lines indicate linear fit and
755 95% confidence intervals, respectively. Color shading indicates local data point density
756 for improved visualization of the extent of overlapping points in the plot, ranging from
757 blue (low density) via green and yellow to red (high density).
758



759

760 **Figure 6. Observed site frequency spectrum of deleterious variants.** The site
761 frequency spectra of (a) probably damaging and (b) loss-of-function variants are shown
762 for the entire exome (white bars) and the MHC region only (blue bars, excluding the
763 classical HLA loci).

764

765