# Extreme distribution of deleterious variation in a historically small and isolated population – insights from the Greenlandic Inuit

Casper-Emil T. Pedersen*, Kirk E. Lohmueller[†‡], Niels Grarup[§], Peter Bjerregaard[**], Torben Hansen[§], Hans R. Siegismund*, Ida Moltke*[1], and Anders Albrechtsen*[1]

*Department of Biology, Section for Computational and RNA Biology, University of Copenhagen, 2200 Copenhagen N, Denmark

[†]Department of Ecology and Evolutionary Biology, University of California, Los Angeles, CA 90095, United States

[‡]Department of Human Genetics, David Geffen School of Medicine, University of California, Los Angeles, CA 90095, United States

[§]The Novo Nordisk Foundation Center for Basic Metabolic Research, Section of Metabolic Genetics, Faculty of Health and Medical Sciences, University of Copenhagen, 2100 Copenhagen E, Denmark

[**]National Institute of Public Health, University of Southern Denmark, 1353 Copenhagen K, Denmark

## Abstract

The genetic consequences of a severe bottleneck on genetic load in humans are widely disputed. Based on exome sequencing of 18 Greenlandic Inuit we show that the Inuit have undergone a severe ~20,000 yearlong bottleneck. This has led to a markedly more extreme distribution of deleterious alleles than seen for any other human population. Compared to populations with much larger population sizes, we see an overall reduction in the number of variable sites, increased numbers of fixed sites, a lower heterozygosity, and increased mean allele frequency as well as more homozygous deleterious genotypes. This means, that the Inuit population is the perfect population to examine the effect of a bottleneck on genetic load. Compared to the European, Asian and African populations, we do not observe a difference in the overall number of derived alleles. In contrast, using proxies for genetic load we find that selection has acted less efficiently in the Inuit, under a recessive model. This fits with our simulations that predict a similar number of derived alleles but a true higher genetic load for the Inuit regardless of the genetic model. Finally, we find that the Inuit population has a great potential for mapping of disease-causing variants that are rare in large populations. In fact, we show that these alleles are more likely to be common, and thus easy to map, in the Inuit than in the Finnish and Latino populations; populations considered highly valuable for mapping studies due to recent bottleneck events.

## Introduction

53

54 Predictions about the consequences of a small population size on genetic variation are

55 among the most fundamental theoretical predictions in population genetics [1–3]. Small

56 populations are more affected by drift than are large populations; therefore, small

57 populations are predicted to carry lower genetic diversity [1]. Additionally, natural selection

58 acting against deleterious alleles is predicted to be less efficient in small populations [2].

59 Finally, the genetic load, defined as the reduction in fitness caused by deleterious variation

60 [4], is predicted to be larger in small populations [1]. With the recent advent of genome-wide

61 sequencing data from humans, it has become possible to test these predictions in human

62 populations and numerous studies have pursued this by comparing the distributions of

63 deleterious alleles across human populations [5–7]. Most of these studies have focused on

64 genetic consequences of the bottlenecks, i.e. rapid large decreases in population sizes

65 that all non-Africans populations went through during the Out-Of-Africa (OOA) dispersal.

66 These studies all agreed that the European population [6–9] and other non-African

67 populations [5] carry lower levels of genetic diversity compared to the African populations

68 [10,11] and thus have a higher genetic load assuming mutations have a recessive effect.

69 However, the studies [7,9,10,12] disagreed in their conclusions about the extent to which

70 genetic load vary across populations assuming mutation have an additive effect and to

71 whether selection is less efficient in small populations.

72 There is no direct way of calculating genetic load from genetic data from natural

73 populations [11]. Instead, summary statistics are used and the conclusions based on

74 different statistics in different studies vary. On one hand, a recent study by Do et al. [8]

75 found no difference between the non-African and African 1000 Genomes populations

76 using the $R_{X/Y}$ statistic, which is monotonically related to the difference in number of

77    derived alleles between the two populations. In line with this finding, two other studies

78    found no significant difference in the number of derived alleles per genome between

79    individuals of European descent and individuals of African American descent [7,8]. Based on

80    this observation, Simons et al. [7] concluded that Europeans and African Americans carry

81    the same amount of genetic load and supported by simulations, they furthermore

82    concluded that genetic load is not affected much by recent population size changes. On

83    the other hand, other studies found slightly, but significantly, more derived deleterious

84    alleles per individual in the European population than in African populations [6,9] implying

85    that the OOA bottleneck has led to a small increase in genetic load in the European

86    population. Similarly, in a recent study of seven populations, Henn et al. [5] reported that the

87    non-African populations on average harbor slightly more derived alleles among alleles with

88    large Genomic Evolutionary Rate Profiling (GERP) scores. Furthermore, Henn et al. [5]

89    reported a small, but significant, difference in load across the seven populations under an

90    additive effect model when load is estimated using selection coefficients approximated

91    based on GERP scores.

92       The consequences of smaller population sizes on the efficacy of selection in

93    humans was first investigated in a study of exome data from 15 African Americans and 20

94    Europeans [13]. That study showed that the proportion of SNPs that are non-synonymous,

95    and thus likely deleterious, is larger in Europeans than in African Americans. This

96    observation was initially interpreted as being, in part, due to less efficient selection during

97    the OOA bottleneck [5] combined with the recent influx of nonsynonymous mutations during

98    the recovery of the bottleneck. Several later studies have come to similar conclusions [5,9,11].

99    In contrast, based on their results about genetic load described above, Do et al. [8]

100   concluded the efficacy of selection was not reduced due to the reduced population size

4

101  during the OOA bottleneck [7,8]. A recent review [11] tried to reconcile the different

102  observations and conclusions by pointing out that different studies have focused on

103  different definitions of the efficacy of selection and use different metrics to quantify it. The

104  review ends by calling for descriptions of empirical patterns of deleterious mutations in

105  other human populations than Europeans and Africans with different population size

106  histories to shed further light on the questions about the effect of small size on selection

107  and load. Motivated by this, we here analyze patterns of deleterious mutations in the

108  Greenlandic Inuit (GI) population based on exome sequencing data from 18 GI individuals.

109      The GI came to Greenland less than 1000 years ago [14,15] and a previous study has

110  shown that their ancestors split from the closest large old-world population, the Han

111  Chinese (CHB) population, some ~23K years before present with only a limited amount of

112  subsequent gene flow [16]. The study also showed that the GI population has a very small

113  effective population size [16]. This suggests that the GI population has likely been small and

114  isolated for a long period of time after the OOA dispersal took place. Analyses of the GI

115  population could therefore be particularly important for resolving the open questions about

116  whether a small effective population size has had any effect on genetic load and selection

117  in human populations.

118      In this paper we first show that the GI population is indeed one of the smallest

119  populations in the world and therefore valuable to analyze. Then we compare different

120  proxies for genetic load between the GI and Europeans, as well as investigate several

121  additional measures related to efficacy of selection. And based on our results we discuss

122  the effect of a small population size on load and selection. Finally, we present analyses

123  that investigate the consequence of our findings to disease genetics, which show that the

124  GI population has great potential for leading to discoveries of new disease related genetic

125  variants that have been missed in large scale GWAS in Europeans and other large

126  populations.

127

## Materials and Methods

### Exome Datasets:

130  We based all of our analyses of the GI population on the high-depth whole exome

131  sequencing data from Moltke et al. [17]. This dataset consists of data from 9 Greenlandic

132  trios, however we restricted our analyses to data from the 18 parents. The GI data was

133  both analyzed alone and together with other datasets: 1) whole exome sequencing data

134  from 18 unrelated Utah residents with Northern and Western European ancestry (CEU), 18

135  Han Chinese individuals from Beijing (CHB) and Yoruban individuals from Ibadan, Nigeria

136  (YRI) from the 1000 Genomes project [18], 2) an exome dataset which includes Mayans,

137  Mbuti and Cambodians from Henn et al. [5] and 3) called genotypes from the 1000 Genome

138  low-coverage data. The analyzed datasets are described in detail below.

139  *GI dataset* This dataset consists of high-depth exome data from 18 GI individuals. To

140  generate it, Moltke et al. [19] performed SNP and indel calling followed by genotype calling

141  for the exome data from the 18 Greenlandic samples using Samtools version 0.1.18 [20].

142  Reads with a mapping quality lower than 30 as well as bases with a base quality lower

143  than 20 had been removed. SNPs had been called using standard settings and genotypes

144  were called based on the highest genotype likelihood. We used VCFtools version 0.1.11 [21]

145  for filtering. We removed SNVs where all individuals were heterozygous, which are likely

146  genotype errors. In addition, we removed sites with sequencing depth lower than 10 for all

147  individuals or higher than 500. The resulting dataset comprise 133.808 SNVs within the

148    exome. This data set was used solely for identifying absolute numbers of SNVs across a

149    number of functional categories and for inferring population size changes though time.

150

151    **Combined exome dataset:**

152    This dataset consists of exome data from both the GI and the CEU population. Further,

153    this dataset also includes exome data from 18 Han Chinese individuals in Beijing (CHB)

154    and 18 Yoruban individuals in Ibadan, Nigeria (YRI). To generate it, we performed joint

155    SNP and genotype calling for the exome data from both populations using ANGSD [22]

156    under the Samtools[18] genotype likelihood model. Reads with a mapping quality lower than

157    30, as well as bases with a base quality lower than 20, were removed. SNVs were called

158    using a likelihood ratio test [22] with a *P*-value cutoff of $10^{-6}$ and genotypes were called

159    based on the highest genotype likelihood. We required a minimum depth of 10 for calling

160    genotypes and removed sites with missing genotypes and sites, which were triploid when

161    including the ancestral allele. This dataset comprises 295.065 SNVs. SFS comparison of

162    1000G genotype calls from the low depth data and our genotype calling from the exome

163    data for the same individuals revealed no differences in proportions (Figure S7). The

164    datasets was polarized using the chimp data available within Seattleseq 138 annotation.

165    For each pairwise comparison only the sites that are polymorphic in the two populations

166    were used.

167

168    **Dataset from 1000 genomes low-depth whole genome sequencing:**

169    This dataset consists of data from five of the 1000 genomes population samples[18]: Finns

170    from Finland (FIN), Peruvians from Lima, Peru (PEL), Gujarati Indians from Houston,

171    Texas (GIH), Utah Residents (CEPH) with Northern and Western Ancestry (CEU), Yoruba

7

172    in Ibadan, Nigeria (YRI) and Han Chinese in Bejing, China (CHB). To generate the

173    dataset, we used VCF files with genotypes calls from the complete phase 3 1000genomes

174    dataset. From this dataset we extracted 11 unadmixed individuals from each of the 5

175    populations. Because of apparent admixture in the Peruvians, we reduced the sample size

176    for each population to 11 individuals. Unadmixed individuals from PEL were selected

177    based on inferred admixture proportions (<5%) using ADMIXTURE. The final dataset

178    comprised 401.821 SNVs within the exome.

179

180    **Exome sequencing data for Mayans, Mbuti and Cambodians:**

181    This dataset consists of data from six samples from each of the populations Mayan, Mbuti

182    and Cambodian. To generate it, we downloaded VCF files with called genotypes for

183    individuals from seven populations recently made publicly available by Henn et al. [5]. From

184    this dataset we extracted data from 6 Mbuti Pygmies, 6 Mayan Indians and 6 Cambodians.

185    More specifically, following Henn et al.[5] we excluded two of the 8 Mayan samples due to

186    admixture, leaving only 6 Mayan samples. To ensure comparability between the SFSs

187    inferred from the different populations, we included only 6 randomly chosen individuals

188    from each of the other populations. The ancestral allele in each dataset was obtained

189    separately using the Pantro2 chimp allele. The final dataset contains genotype data for

190    194.278 SNVs.

191

192    **Annotation of variants:**

193    We divided variant sites into four categories based on the functional category annotated to

194    the derived allele: synonymous, a combined category of non-coding exons (NCE) and the

195    leading and trailing untranslated region (UTR), missense and loss of function (LoF).

196    Variant sites which belong to other functional annotations were excluded from further

197    analysis. We assigned a variant to the putatively most deleterious annotation category

198    when there were multiple splice variations. We also included an additional category for

199    loss of function where all isoforms are annotated to be loss of function (denoted LoF[a]),

200    meaning that these sites are more likely to be deleterious than the sites in the LoF

201    category. The five categories thus range from mutations that are expected to be neutrally

202    evolving to mutations that are expected to be highly deleterious.

203    We also divided variants in categories using GERP scores [23]. These GERP score

204    measures conservation across a phylogeny of 35 mammalian species excluding humans [5].

205    Specifically, GERP scores represent the deficiency in numbers of substitutions in

206    functional loci compared to that of the number of substitutions seen in neutral DNA. This

207    discrepancy is then regarded as a sign of functional constraint and thus lower degrees of

208    substitution saturation across the phylogeny will reflect higher levels of purifying selection

209    [5,24]. We used the same approach as Henn et al. [10]: we retrieved the GERP scores from the

210    UCSC browser and grouped variants into four categories according to their GERP score,

211    and thus how deleterious they are predicted to be. The categories are "neutral" (GERP <

212    2), "moderate" ($2 \leq$ GERP < 4), "large" ($4 \leq$ GERP < 6) and finally "extreme" (GERP $\geq$ 6).

213

214    **Population size inference:**

215    We performed inference of population size over time on the GI dataset using the site

216    frequency spectrum (SFS) based method called stairway plot [25]. We performed the analyis

217    using default settings.

218

219

9

220 **Genetic load approximating summary statistics the number of derived alleles:**

221 We used two different proxies for genetic load based on counts of derived alleles: the total

222 number of derived alleles per individual, which has previously been used to quantify load

223 under an additive model [9] and the number of homozygous-derived genotypes per

224 individual, which can equivalently be used to quantify load under a recessive model.

225

226 **The GERP score load:**

227 We also used what we will denote GERP score load. To calculate this, we used the

228 GERP-based grouping of deleterious sites described above and translated the three non-

229 neutral GERP score categories into values of selection coefficients as suggested by Henn

230 et al. [10]. Specifically, the extreme category was assigned $s = 1 \times 10^{-2}$, the large category

231 was assigned $s = 4.5 \times 10^{-3}$, while the moderate category was assigned $s = 4.5 \times 10^{-4}$. These

232 assignments were then used in the equation from Kimura[1]

$$Load = 1 - w = 1 - (1 - 2q(1 - q)sh - sq^2)$$

234 which is based on the underlying model where the fitness of each genotype is determined

235 as 1, 1-$hs$ and 1-$s$ for AA, Aa and aa, respectively. We used this equation assuming two

236 models each with an extreme level of dominance. In the first model $h$ was set to 0.5, i.e.

237 we assumed alleles to act additively. In the other model $h$ was set to 0, i.e. we assumed

238 deleterious alleles to be completely recessive. To estimate standard error for our GERP

239 score loads estimates we used a weighted uneven block jackknife, with block sizes of 5Mb

240 to correct for correlation among neighboring sites [26].

241

242 $\pi_{var}$ **calculation:**

10

243    We calculated the nucleotide diversity for the variable sites, $\pi_{var}$, for all the populations

244    included in the SFS comparison in Figure 2. We used the following equation for our

245    calculations:

$$\pi_{var} = \frac{1}{\binom{2n}{2}} \sum_{i=1}^{2n} x_i \, (2n - i) i$$

246    where $x_i$ is the proportion of sites in the ith category of the site frequency spectrum and *n*

247    the number of individuals included in the analysis. We calculated $\pi_{var}$ rather than $\pi$ for all

248    sites, because it is not possible to calculate $\pi$ for all sites from the 1000 genomes

249    genotype calls since there is not information for the invariables sites.

250

251

**Simulations:**

252

253  We performed simulations using the forward-time simulation software Selection on Linked

254  Mutations (SLIM) [27]. We used the population size and the demographic history estimated

255  for the GI and the CHB by Fumagalli et al. [16]. We initiated simulations with a haploid

256  population size of 9395 at 29459 generations before the sampling time (present time),

257  allowing sufficient time for saturation of neutral mutations. During the OOA bottleneck

258  population size was reduced to 5443 haploid genotypes. The CHB population did not

259  change population size up until present time. To simulate the additional bottleneck, which

260  the GI population underwent after splitting from CHB, we introduced an additional

261  reduction in population size 930 generations before the present, resulting in a population

262  size of 1550 haploid individuals. We sampled 36 haploid genotypes from each of the

263  simulated GI and CHB populations, allowing us to capture comparable levels of genetic

264  diversity as the real data where we have 18 diploid samples from each. For this simulated

265  demographic scenario, we chose to simulate many sites under selection: 90% were

266  simulated as neutral mutations and 10% as being under selection. For the deleterious

267  mutations, we varied the dominance coefficient ($h$) and the selection coefficient ($s$) to

268  mimic the effects of three modes of inheritance: additive ($h$=0.5), near-recessive ($h$=0.1)

269  and recessive ($h$=0). Furthermore, for each of these modes, we varied the selection

270  coefficient, $s$, across four levels, reflecting increasingly harmful functional effects

271  ($s$=0.0002; $s$=0.002; $s$=0.02; $s$=0.2). Note, that because SLIM have a different underlying

272  model of selection implemented than the one we mentioned above, hence the selection

273  coefficients provided here.

274   Each of the scenarios, i.e. combination of $h$ and $s$ value, was simulated 10000 times on a

275  sequence of 100k sites using a recombination rate of $1.2 \times 10^{-8}$ bp. and mutation rates of

12

276  1.38 × 10$^{-8}$ per generation. Both numbers of derived alleles and numbers of homozygous-

277  derived genotypes are reported excluding mutations that were fixed in both populations.

278  Furthermore, we simulated a scenario where all mutations were set to neutral in GI after

279  the split between GI and CHB. For this simulation, we sampled 36 haploid genotypes both

280  in the pooled population before the split and in each of the GI and CHB populations.

281

282  **Allele frequencies in GI compared to Europeans, East Asians, Finns and Latinos:**

283  To compare allele frequencies for shared alleles across populations, we used the

284  information available for more than 9.3 million sites based on exome data in the ExAC

285  browser. This dataset has sufficient number of samples for European and East Asian

286  populations to allow reliable calculation of rare allele frequencies.

287

288  ## Results

289  **Inference of demography and SFS:**

290  First we inferred population size of the GI population over time from exome data for 18 GI

291  using the stairway plot method [25]. The results suggest that the GI population experienced

292  a marked decrease in size some 22-24 Kya (Figure 1), remained small for more than

293  20,000 year and only recently started to increase in size. We note that there is a

294  considerable amount of uncertainty for recent population size estimates, which makes it

295  difficult to assess the extent of the recent increase.

296  Next, we compared the Site Frequency Spectrum (SFS) of the GI population to the

297  SFSs of a number of other human populations from the 1000 genomes project (Figure 2)

298  and Cambodians and Mayans from Henn et al. [5] (Figure S8). Because the SFS is affected

299  by demographic history [28,29], this comparison should reveal the extent to which the recent

300    very long bottleneck has made the GI more extreme than other OOA populations. These

301    comparisons show that the SFS for the GI population is flatter than for all the other

302    examined populations, which is also reflected in markedly higher nucleotide diversity for

303    the variable sites ($\pi_{var}$), which is inversely correlated with $\pi$ for all sites. Thus compared to

304    other populations, the GI population has had both a larger depletion of rare variants, as

305    well as a higher increase in allele frequencies for remaining variants. Hence the GI

306    population shows evidence for having undergone a more extreme population bottleneck

307    than the other populations, which means that analyses of this population are likely to be

308    valuable in resolving the open questions about whether small effective population sizes

309    have had an effect on genetic load and selection in human populations.

310

**311    Comparison of genetic load in GI and CEU:**

312    Motivated by these initial findings, we compared the genetic load between GI and CEU to

313    investigate to which extent, if any, the long history with a reduced population size has led

314    to an increased genetic load in the GI population. For this comparison we used the exome

315    data from the 18 GI combined with exome data from 18 individuals with European ancestry

316    (CEU). Since there is no direct way to calculate load from genetic data, we examined

317    several different statistics, most of which have recently been proposed as approximations

318    to genetic load.

319        First, we examined load statistics that assume that all alleles have an additive

320    effect. In particular, we looked at the average number of derived alleles per individual,

321    which has previously been used both by Fu et al. [9] and Do et al. [8]. We performed the

322    comparison between GI and CEU by calculating the ratio of this statistic in GI and CEU.

323    Here a ratio above one would indicate that the average number of derived alleles in GI is

14

324  higher than that in CEU. We calculated this ratio both for all sites and for sites in 9 different

325  subcategories that represent a range of different levels putative deleterious effect: five

326  categories based on functional annotation and four based on GERP scores (for details see

327  Material and Methods). We did not find a significant difference between the numbers of

328  derived alleles per individual between the two populations for any of the different functional

329  annotation-based categories, i.e. the ratios did not differ significantly from 1 (Table 1 and

330  Figure 3, Figure S4, Figure S1, S2, Table S1 and Table S2). The ratio of derived alleles in

331  the more deleterious GERP categories show a 1-4% increase in GI compared to CEU

332  (Figure 4), but this difference is not significant. Similar results are obtained from the $R_{X/Y}$

333  ratio from Do et al. [8] Hence, the results from both these statistics suggest that, under an

334  additive effect model, the load is the same in the across populations (Table 1, Table S1

335  and Table S2). However, we note that the percentage of the derived alleles per individual

336  that come from sites fixed for the derived allele differs markedly between the two

337  populations. This is seen for all sites as well as in most functional subcategories (e.g. LoF

338  alleles and GERP scores higher than 6; Table 1; column 4). We also note that the average

339  number of derived alleles per SNV differs significantly between the two populations in all

340  categories (Table 1; column 6).  All the ratios are significantly larger than 1, indicating that

341  per SNV, the derived allele is at higher frequency in GI compared to CEU. This

342  observation combined with our previous observation that the average number of derived

343  alleles per individual is the same in the two populations, reveals that the derived alleles in

344  GI are from fewer sites with higher derived allele frequencies and thus that the load under

345  an additive model is distributed differently in the two populations.

346      Next, we aimed to compare the genetic load under the assumption that deleterious

347  alleles have a recessive effect. To do this we examined the average number of

348    homozygous derived genotypes per individual. The ratios of this statistic were significantly

349    higher than 1 (Table 1; column 2 and Figure 3, Figure S4), indicating an accumulation of

350    homozygous derived genotypes in the GI population. This was seen for all categories,

351    even for the LoF categories. In fact, the ratio is significantly higher in sites with large

352    deleterious effects, i.e. GERP scores above 4, compared to neutral sites, i.e. sites with

353    GERP scores below 2 (see asterisk in Figure 4). These observations suggest that the load

354    is increased in GI, if the mode of selection is recessive.

355        We also looked at the load statistic proposed by Henn et al. [10], which we denote as

356    GERP score load. This statistic approximates load using the original definition by Kimura

357    et al. [1], where each GERP score category is translated into a selection coefficient. We

358    coupled these groups of selection coefficients with two models: one where $h$ was set to

359    0.5, i.e. assuming alleles have an additive effect and one where $h$ was set to 0; i.e.

360    assuming alleles have a recessive effect (for details see Materials and Methods). When

361    doing so, we saw no significant difference between the two populations under the additive

362    model (10.1 (S.E. 0.4) for CEU vs 10.2 (S.E. 0.4) for GI), but a 22% increase in the GERP

363    score load under the completely recessive model (5.1 (S.E. 0.2) vs 6.2 (S.E. 0.3)) (Figure

364    3). Furthermore, putatively deleterious SNVs (GERP score > 2) that contribute to the

365    genetic load are much more common in the GI population compared to the CEU

366    population (Figure 3). Hence, the results obtained based on GERP scores lead to

367    qualitatively the same conclusion as the other statistics.

368        We note that we also performed similar load comparisons to 18 individuals with

369    East Asian ancestry (CHB) and the 18 individuals with African ancestry (YRI) to ensure

370    that our conclusions were not artifacts of particular features of the CEU. In particular, we

371    produced the equivalents of Table 1 and figure 4 for these two populations (Table S1,

16

372   Table S2, Figure S1 and Figure S2). When doing so we reached the same conclusions for

373   both populations with one exception: when comparing GI and CHB, the ratio of

374   homozygous-derived genotype counts is indeed higher in sites with large deleterious

375   effects, i.e. GERP scores above 4, compared to neutral sites, but not significantly so

376   (Figure S1). However, for YRI this difference is significant (asterisk in Figure S2).

377

378   **Comparison of genetic load based on simulations:**

379   To investigate these results further, and try to assess what can be concluded from them,

380   we performed simulations of GI and CHB populations using the demographic history of

381   these two populations inferred in Fumagalli et al. [16] We performed the simulations of 12

382   different selection scenarios that varied both in the effect model used, i.e. $h$, and the

383   selection coefficients, s, that reflect how deleterious the alleles are. Based on the

384   simulated data we calculated the ratio of the number of derived alleles between the two

385   populations and this ratio was not significantly different from 1 in any of the scenarios

386   (Table S6).

387        We also counted the number of homozygous derived genotypes and, as expected,

388   we observed a significantly higher number GI than in CHB (Figure S8-9). The increase is

389   about 19% even without selection (Table S6).

390        Finally, we calculated the true load based on the true effect model and the true

391   selection coefficient for each site. Interestingly the ratio of this true load in GI and CHB is

392   significantly higher than 1 in all scenarios, indicating a higher load in GI than in CHB. This

393   difference decreases as $h$ increases, however, importantly, even in the scenarios where

394   additive selection ($h$=0.5) was simulated, the ratio of loads were all significantly above 1,

395   with the increase varying from 1 to 4% (Figure S9 and Table S6) depending on the

17

396  strength of selection simulated. The highest load difference was observed for the scenario

397  with $s$=0.002.

398  These results are interesting because they suggest that load statistics, like the

399  number of derived alleles, will not necessarily reveal if there is a difference in load

400  between two populations. Thus the results suggest that the fact that GI do not have higher

401  additive load statistic than CEU in our data cannot necessarily be used to conclude that

402  there is no difference in load even if all alleles have an additive effect. For the same

403  reason, it cannot necessarily be used to conclude that selection has not been acting less

404  efficiently in the GI.

405

406  **Effectiveness of selection:**

407  To investigate to what extent the long recent bottleneck in GI has affected the

408  effectiveness of selection, we looked at two other summary statistics, which have both

409  previously been used to address this question in the context of the OOA bottleneck: SFSs

410  and the proportion of non-synonymous to synonymous mutations[13].

411  We first compared SFSs for GI and the CEU. Specifically we made SFSs for each

412  of five the functional categories also used in our previous analyses: synonymous, a

413  combined category of non-coding exons (NCE) and the leading and trailing untranslated

414  region (UTR), missense and finally two different LoF categories (Figure S3 and Figure S4).

415  These SFSs show that GI have a larger proportion of alleles found in higher frequencies,

416  not only overall, but also within all of the functional categories including LoF, and thus

417  potentially highly deleterious alleles. And consistent with the results described previously,

418  the SFSs also show that the GI population has a higher proportion of fixed putatively

419  deleterious derived alleles in all the functional categories. In contrast, the CEU population

18

420    has a clear skew towards rare variants. This difference, may reflect the effects of less

421    efficient purifying selection in GI and/or recent growth in CEU [6,29]. However, we note that

422    based on simulations similar to the ones presented above, but where each deleterious

423    allele is set to neutral (changing the selection coefficient to 0) in GI after the start of the GI

424    bottleneck (Figure S10), we observed a fairly similar SFS to the one where selection

425    continues to happen after the bottleneck. Thus based on the simulations we do not

426    necessarily expect to be able to observe if selection is less effective based on large

427    differences in the SFS of the different functional categories. Next, we compared the ratio of

428    non-synonymous to synonymous SNPs in the GI populations to that of the CEU population

429    using the Combined GI and CEU dataset, following what was done by Lohmueller et al. [13]

430    (Table S4). Specifically, we tested for a difference in this ratio between the two populations

431    both among all sites and among sites from different GERP scores categories. We found

432    that ratio of the number of non-synonymous to synonymous SNPs is significantly higher in

433    CEU than in GI, when considering all sites (first row in Table S4). When considering

434    different GERP score categories, the non-synonymous to synonymous ratio in GI becomes

435    even lower relative to that seen in CEU with increasing levels of deleteriousness.

436       We further investigated this pattern by defining neutral sites as those with GERP

437    scores below 2 and deleterious variants as those with GERP scores above 4 and then

438    calculating the ratio of deleterious to neutral SNPs for both GI and CEU. When doing so

439    the deleterious to neutral SNP ratio is significantly lower in the GI compared to the CEU

440    (Table S4; row 2) and is almost as low as the ratio of non-synonymous to synonymous

441    sites among the most deleterious GERP scores (Table S4; row 5 and 6). Hence, there

442    seem to be a higher proportion of deleterious variants in CEU than in GI, particularly for

443    highly deleterious variants. This is compatible with the notion that selection has not been

19

444    acting less efficiently in GI than in CEU. Hence the different statistics point towards

445    different conclusions, which we will discuss later.

446

447    **Consequences for disease mapping:**

448    Regardless of whether selection has acted less efficiently in GI or not, our results clearly

449    show that individual deleterious variants tend to have higher frequency in GI compared to

450    CEU (Figure 3). This is potentially of high importance for quantitative trait mapping and

451    disease mapping as it means that mapping in the GI population will in some cases be

452    more powerful in the GI population than the CEU population [30] just like it is well-known to

453    be the case for other historically isolated populations like the Finnish population and the

454    Native Americans. To investigate the potential of GI, and compare it to that of other

455    populations, we divided alleles into bins according to their frequency in larger reference

456    populations and for each bin determined how often the alleles are common among GI

457    compared to how often they common in other populations. To do this we used allele

458    frequencies from the Exome aggregation consortium (ExAC). As can be seen in Figure 5

459    and Figure S9-S11, we find that rare alleles among Non-Finnish Europeans are more likely

460    to be common in GI than in Finland, East Asia and Latinos. This pattern is especially

461    pronounced for variants that are extremely rare in Non-Finnish Europeans (MAF $< 1 \times 10^{-3}$).

462    Similar patterns are seen when looking at variants that are rare in East Asians. This

463    suggests that the GI population does indeed have potential to provide increased power to

464    detect alleles that are rare in Non-Finnish Europeans and East Asians. And it suggests

465    that this potential is even bigger than that of both Latinos and Finns; populations that are

466    considered particularly powerful for mapping.

467  To pursue this potential of the GI population in a simple manner we investigated all the

468  LoF including frameshift variants that are present in at least 2 copies in the GI dataset and

469  are rare, but not absent, in both the European and EAS populations (MAF below 0.5%)

470  (Table 2, Table S5). We did this because such alleles are of potential interest for the

471  Europeans and EAS, but difficult to map and investigate in these populations. This

472  approach led to the identification of 6 SNVs and 14 indels. One of the SNVs, located in

473  *TBC1D4*, was recently shown to have a very large impact on type 2 diabetes[15] and one of

474  the indels, located in *SI,* was recently shown to have large impact on a sucrase-isomaltase

475  deficiency [31]. Other interesting variants include a SNV in *SEMA4C* which gene is

476  associated with neonatal lethality[17,32] and SNVs in *CRYGA* and *USP45* which might be

477  involved in cataract and DNA repair, respectively [33,34].

478

## Discussion

480  We have performed analyses of exome data from the GI population with the aim of

481  investigating what consequences small population size has had on genetic variation in

482  humans. We first investigated the demographic history of the GI population and found that

483  it dramatically decreased in size ~23.000 years ago and remained small for more than 20

484  thousand years, which corroborates previously reported results [16]. Furthermore, we

485  observed a flattening of the SFS for populations with increasing distance to Africa, and that

486  the SFS for the GI population was even more extreme than those for the equally distant

487  Mayan and Peruvian populations. This suggests that the GI population underwent a strong

488  bottleneck after it split from the Native Americans. This is consistent with our estimates of

489  the GI population size over time (Figure 1), where the bottleneck persisted long after the

490   split from Native Americans. This finding makes the GI population highly relevant to the

491   study of how demography has shaped deleterious variation in humans.

492        Motivated by this we investigated whether the severe bottleneck has affected the

493   genetic load of the GI population by comparing it to CEU. In this context, it is important to

494   re-iterate that genetic load cannot be directly calculated from exome data [11]. It can only be

495   approximated through predictions of whether and to what extent alleles are deleterious,

496   e.g. via functional annotation. Also, the approximations depend on the underlying effect

497   model, i.e. whether the effect of the deleterious alleles is additive or recessive. We

498   therefore compared the genetic load between the GI and CEU populations using several

499   different proxy statistics. Our analyses reveal no significant differences between the two

500   populations when using proxy statistics that assume an additive effect, like the number of

501   derived alleles, but they do show large significant differences when using proxy statistics

502   for load that assume a recessive effect, like the number of homozygous-derived genotypes

503   (Table 1). Importantly, we observed similar results when comparing GI to CHB and YRI

504   (Table S1 and Table S2). These observations are in line with the results from several

505   previous studies [7,8] and are consistent with our simulations where we observe no

506   significant effect on the number of derived alleles and a large significant difference in

507   number of homozygous derived genotypes (Figure 4, Figure S1, Figure S2, Figure S9 and

508   Table S6). When looking at sites with GERP scores or functional annotations, which

509   suggest that that they harbor highly deleterious alleles we observed a larger increase in

510   number of derived alleles in GI than in CEU, although the increase is small and non-

511   significant using jackknife resampling. Therefore, our results are not necessarily conflicting

512   with Henn et al. [5] that found that the number of deleterious alleles increases with distance

513   from Africa. However, we would have expected to see a stronger difference in GI than in

514    the Mayan used by Henn et al. [5]. Though, it is possible that we would see such evidence

515    in our ratios if we analyzed a larger number of samples.

516         Interestingly, our simulations show a significant increase in the true load in the

517    simulated GI individuals compared to the simulated CHB individuals. Even under an

518    additive model, there is a small but significant increase of 1-4% depending on the selection

519    coefficient. This suggests that small population size over a long period of time can indeed

520    lead to increased load even if all alleles have an additive effect. These simulation results

521    are in conflict with simulation-based conclusions from Simons et al. [7], based on deleterious

522    alleles counts. However, our results fit very well with recent results from Harris and Nielsen

523    [35]: when simulating Neanderthals and humans and assuming a severe bottleneck in

524    Neanderthal that lasted many times longer than the bottleneck in GI, they observed at

525    large significant decrease in fitness among Neanderthals even under an additive effect

526    model.  Importantly, our simulation results also revealed that the significant increase in

527    load observed under additive selection, did not lead to a significant difference in the

528    number of derived alleles. This suggests that the fact that we do not see a significant

529    increase in the number of derived alleles in GI compared to CEU (or the other two

530    populations) does not necessarily mean that there is no difference in load under an

531    additive model. Furthermore, it supports the conclusion from Lohmueller [11] that part of the

532    reason why different studies reach different conclusions about load is that they are based

533    on different approximations of load – some of which may be far from ideal.

534         Either way, our results combined with previous results suggest, that if there are

535    increases in genetic load under an additive model due to the small population size and

536    bottlenecks, these increases are small even for an extreme population like the GI.  Our

537    results also clearly show that the extent to which the load is increased in GI due to the

538     severe bottleneck depends heavily on the true underlying effect model. Hence to fully

539     answer this question more knowledge about the true effect model is needed.

540         The above conclusions make it difficult to use comparisons of our overall load

541     estimates to make conclusions about whether or not the severe bottleneck in GI has led to

542     less efficient selection. We therefore also looked at other statistics to address that

543     question. One of these was the ratio of non-synonymous to synonymous SNPs. If

544     selection has acted less effectively in GI compared to larger populations like CEU then we

545     would expect this ratio to be higher in the GI population than in CEU. We observe the

546     opposite, which could suggest that selection has not acted less efficiently in GI. However,

547     we note that the ratio of non-synonymous to synonymous SNPs is highly sensitive to

548     demographic changes and that our observations could also be explained by the recent

549     explosive population growth in CEU [11]. Importantly, we also note that if selection has acted

550     equally efficiently in the two populations, then we would expect that the ratio of

551     homozygous-derived genotypes in the two populations would be the same across

552     categories with different levels of deleterious effect. However, we find that this is not the

553     case. On the contrary, the ratio of homozygous-derived genotypes is significantly higher in

554     sites with large deleterious effects, i.e. GERP scores above 4 and below 6 as compared to

555     neutral sites, i.e. sites with GERP scores below 2 (see asterisk in Figure 4). This result

556     was also seen when comparing GI and YRI (asterisk in Figure S2) but not when

557     comparing GI and CHB (see Figure S1). In the latter case there is an increase in the ratio

558     between the two categories but it is not significant (Figure S1). This suggests that the

559     impact of selection is not the same in GI compared to the YRI and CEU populations and

560     potentially not CHB either. More specifically, it suggests that selection against deleterious

561     alleles has acted less efficiently in GI under a recessive model compared to the other

562    populations. This argument is supported by the fact that we observed a similar pattern in

563    our simulations. Here the increase in the homozygous-derived genotype ratio for the sites

564    under selection (denoted as load ratio in Supplementary Table 6) under a recessive model

565    is particularly large for alleles with selection coefficient $s$=0.002, which correspond to a

566    GERP score of 4-6 [5](Figure S9-S11). This load difference would also explain the higher

567    allele frequencies of deleterious alleles frequencies in GI compared to in CEU (Figure 3).

568    Hence these results combined suggest that selection may indeed have acted less

569    efficiently, and that the higher ratio of synonymous to non-synonymous sites GI is most

570    likely explained by the population growth in CEU as has been suggested in other studies

571    as well [12].

572         Load is a somewhat abstract concept, none-the-less it has very concrete

573    connections to disease and although our comparisons of load in GI and CEU may not be

574    entirely conclusive, they have clear consequences in the context of disease. Our load

575    comparisons revealed that a lot of variants are lost, also deleterious variants, and that the

576    deleterious alleles that are not lost are in higher frequency in the GI population than in the

577    CEU population. This has several disease related implications for the Greenlandic

578    population. First, it has implication for disease risk. For complex disorders, where many

579    loci are involved, we do not expect a large difference in genetic load and thus disease risk.

580    This can be illustrated by the *TBC1D4* variant's impact on the type 2 diabetes prevalence.

581    Despite its large effect and relatively high frequency [17] (23% MAF in GI) the Greenlandic

582    population has not had a historically high incidence of type 2 diabetes in Greenland [36]. We

583    suspect that the absence of other type 2 diabetes variants in GI are compensating for the

584    presence of the common variants, but acknowledge that other variants may be found if we

585    investigated a larger amount of data. However, for the rarer more monogenetic traits, we

586  suspect that the genetic load, and thus disease risk, will be very different in Greenland.

587  Because fewer variants are involved, the variance of the prevalence of the trait will be

588  much greater. This means that such diseases will likely either entirely absent from the

589  population or be more prevalent. The SI frameshift variations effect on sucrase-isomaltase

590  deficiency is an example of the latter. SI deficiency is a very rare disorder in large

591  populations, but is estimated to affect 5-10% of individuals in Greenland [37,38]. This is

592  presumably solely due to a single SI frameshift variant, which we find in a homozygous

593  state in 3 of the individuals among the 9 GI trios (Table S5).

594       It also has implication for our ability to characterize the function of certain genes,

595  including disease related genes. Due to the higher frequency of deleterious alleles, even

596  LoF alleles, the GI are enriched for homozygous functional knockouts compared to the

597  CEU. This enables investigation of the function of the genes that harbor such mutations.

598       Finally, it has implication for disease mapping. It has long been acknowledged that

599  populations, which have undergone a recent population bottleneck and therefore, like the

600  Greenlandic population, carry deleterious alleles with higher allele frequency compared to

601  larger populations are useful in disease studies, because the increased allele frequency

602  (along with increased LD) leads to increased power in association testing. The Finnish

603  population is one such population, where two novel mutations have been found to be

604  associated with lipoprotein levels [39]. Other examples included various Native American

605  populations where studies have also led to the detection of several novel disease variants.

606  Interestingly, our analyses show that GI outperform both Finns and Latinos in terms of the

607  chance of providing improved power in association testing due to higher allele frequencies

608  (Figure 5, Figure S11-S14). Thus likely, studies aiming to find novel association will have

26

609   even better power to do so in the GI population compared to previously studied isolated

610   populations.

611         Our simple screening of the LoF variants provides a simple demonstration how

612   useful GI can be for disease mapping. It also clearly shows that studies of the Greenlandic

613   population can be used to identify alleles of large effect. However, populations like the

614   Greenlandic are also useful for identifying variants with lower effect sizes, because, as this

615   study shows, such alleles are much more likely to be of high frequency than are alleles of

616   large effect. Importantly, since our analyses of allele frequencies among GI compared to

617   Finns and Latinos were carried out for alleles that are rare, but indeed present among

618   Europeans and East Asians our results also suggests that disease mapping in the

619   Greenlandic population has great potential to lead to identification of variants that are also

620   important in larger populations, like the Europeans and East Asians. Thus, all in all our

621   results show that studies of GI population constitute a particularly promising approach in

622   future disease mapping.

623

624

625

626

627

628

629

## Supplemental Material

Our supplemental data comprise 14 figures and 6 tables.

## Acknowledgements

## Web Resources

Henn et al.[5] https://ecoevo.stonybrook.edu/hennlab/data-software/

ExAC browser

ftp://ftp.broadinstitute.org/pub/ExAC_release/release0.3/ExAC.r0.3.sites.vep.vcf.gz

GERP scores http://hgdownload.cse.ucsc.edu/gbdb/hg19/bbi/All_hg19_RS.bw

## References

1. Kimura, M., Maruyama, T., and Crow, J.F. (1963). Mutation load in small populations. Genetics *48*, 1303–1312.

2. Ohta, T. (1973). Slightly deleterious mutant substitutions in evolution. Nature *246*, 96–98.

3. Charlesworth, B. (2009). Fundamental concepts in genetics: effective population size and patterns of molecular evolution and variation. Nat. Rev. Genet. *10*, 195–205.

4. Morton, N.E., Crow, J.F., and Muller, H.J. (1956). An estimate of the mutational damage in man from data on consanguineous marriages. Proc Natl Acad Sci U S A *42*, 855–863.

5. Henn, B.M., Botigué, L.R., Peischl, S., Dupanloup, I., Lipatov, M., Maples, B.K., Martin, A.R., Musharoff, S., Cann, H., Snyder, M.P., et al. (2015). Distance from sub-Saharan Africa predicts mutational load in diverse human genomes. Proc. Natl. Acad. Sci. *113*, E440–E444.

6. Boyko, A.R., Williamson, S.H., Indap, A.R., Degenhardt, J.D., Hernandez, R.D., Lohmueller, K.E., Adams, M.D., Schmidt, S., Sninsky, J.J., Sunyaev, S.R., et al. (2008). Assessing the evolutionary impact of amino acid mutations in the human genome. PLoS Genet. *4*, e1000083.

7. Simons, Y.B., Turchin, M.C., Pritchard, J.K., and Sella, G. (2014). The deleterious mutation load is insensitive to recent population history. Nat. Genet. *46*, 220–224.

8. Do, R., Balick, D., Li, H., Adzhubei, I., Sunyaev, S., and Reich, D. (2015). No evidence

663  that selection has been less effective at removing deleterious mutations in Europeans than
664  in Africans. Nat. Genet. *47*, 126–131.
665  9. Fu, W., Gittelman, R.M., Bamshad, M.J., and Akey, J.M. (2014). Characteristics of
666  neutral and deleterious protein-coding variation among individuals and populations. Am. J.
667  Hum. Genet. *95*, 421–436.
668  10. Henn, B.M., Botigué, L.R., Bustamante, C.D., Clark, A.G., and Gravel, S. (2015).
669  Estimating the mutation load in human genomes. Nat. Rev. Genet. *16*, 333–343.
670  11. Lohmueller, K.E. (2014). The distribution of deleterious genetic variation in human
671  populations. Curr. Opin. Genet. Dev. *29*, 139–146.
672  12. Lohmueller, K.E. (2014). The Impact of population demography and selection on the
673  genetic architecture of complex traits. PLoS Genet. *10*, e1004379.
674  13. Lohmueller, K.E., Indap, A.R., Schmidt, S., Boyko, A.R., Hernandez, R.D., Hubisz,
675  M.J., Sninsky, J.J., White, T.J., Sunyaev, S.R., Nielsen, R., et al. (2008). Proportionally
676  more deleterious genetic variation in European than in African populations. Nature *451*,
677  994–997.
678  14. Raghavan, M., DeGiorgio, M., Albrechtsen, A., Moltke, I., Skoglund, P., Korneliussen,
679  T.S., Gronnøw, B., Appelt, M., Gulløv, H.C., Friesen, M., et al. (2014). The genetic
680  prehistory of the new world arctic. Science *345*, 1255832.
681  15. Gulløv, H.C. Grønlands forhistorie (Gyldendal, Copenhagen, 2004).
682  16. Fumagalli, M., Moltke, I., Grarup, N., Racimo, F., Bjerregaard, P., Jørgensen, M.E.,
683  Korneliussen, T.S., Gerbault, P., Skotte, L., Linneberg, A., et al. (2015). Greenlandic Inuit
684  show genetic signatures of diet and climate adaptation. Science *349*, 1343–1347.
685  17. Moltke, I., Grarup, N., Jørgensen, M.E., Bjerregaard, P., Treebak, J.T., Fumagalli, M.,
686  Korneliussen, T.S., Andersen, M.A., Nielsen, T.S., Krarup, N.T., et al. (2014). A common
687  Greenlandic TBC1D4 variant confers muscle insulin resistance and type 2 diabetes.
688  Nature *512*, 190–193.
689  18. McVean, G.A., Altshuler (Co-Chair), D.M., Durbin (Co-Chair), R.M., Abecasis, G.R.,
690  Bentley, D.R., Chakravarti, A., Clark, A.G., Donnelly, P., Eichler, E.E., Flicek, P., et al.
691  (2012). An integrated map of genetic variation from 1,092 human genomes. Nature *491*,
692  56–65.
693  19. Moltke, I., Fumagalli, M., Korneliussen, T.S., Crawford, J.E., Bjerregaard, P.,
694  Jørgensen, M.E., Grarup, N., Gulløv, H.C., Linneberg, A., Pedersen, O., et al. (2015).
695  Uncovering the genetic history of the present-day greenlandic population. Am. J. Hum.
696  Genet. *96*, 54–69.
697  20. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,
698  Abecasis, G., Durbin, R., and Subgroup, 1000 Genome Project Data Processing (2009).
699  The sequence alignment/map format and SAMtools. Bioinformatics *25*, 2078–2079.
700  21. Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A.,
701  Handsaker, R.E., Lunter, G., Marth, G.T., Sherry, S.T., et al. (2011). The variant call
702  format and VCFtools. Bioinformatics *27*, 2156–2158.
703  22. Korneliussen, T., Albrechtsen, A., and Nielsen, R. (2014). ANGSD: Analysis of Next
704  Generation Sequencing Data. BMC Bioinformatics *15*, 356.
705  23. Cooper, G.M., Stone, E.A., Asimenos, G., Green, E.D., Batzoglou, S., and Sidow, A.
706  (2005). Distribution and intensity of constraint in mammalian genomic sequence. Genome
707  Res. *15*, 901–913.
708  24. Goode, D.L., Cooper, G.M., Schmutz, J., Dickson, M., Gonzales, E., Tsai, M., Karra,
709  K., Davydov, E., Batzoglou, S., Myers, R.M., et al. (2010). Evolutionary constraint
710  facilitates interpretation of genetic variation in resequenced human genomes. Genome

711  Res. *20*, 301–310.
712  25. Liu, X., and Fu, Y.-X. (2015). Exploring population size changes using SNP frequency
713  spectra. Nat. Genet. *47*, 555–559.
714  26. Busing, F.M.T.A., Meijer, E., and Leeden, R.V.D. (1999). Delete- m Jackknife for
715  Unequal m. Stat. Comput. *9*, 3–8.
716  27. Messer, P.W. (2013). SLiM: Simulating Evolution with Selection and Linkage. Genetics
717  *194*, 1037–1039.
718  28. Marth, G.T., Czabarka, E., Murvai, J., and Sherry, S.T. (2004). The Allele Frequency
719  Spectrum in Genome-Wide Human Variation Data Reveals Signals of Differential
720  Demographic History in Three Large World Populations. Genetics *166*, 351–372.
721  29. Nielsen, R., Bustamante, C., Clark, A.G., Glanowski, S., Sackton, T.B., Hubisz, M.J.,
722  Fledel-Alon, A., Tanenbaum, D.M., Civello, D., White, T.J., et al. (2005). A scan for
723  positively selected genes in the genomes of humans and chimpanzees. PLoS Biol. *3*,
724  e170.
725  30. Hong, E.P., and Park, J.W. (2012). Sample Size and Statistical Power Calculation in
726  Genetic Association Studies. Genomics Inform. *10*, 117.
727  31. Marcadier, J.L., Boland, M., Scott, C.R., Issa, K., Wu, Z., Mcintyre, A.D., Hegele, R.A.,
728  Geraghty, M.T., and Lines, M.A. (2014). Congenital sucrase—isomaltase deficiency:
729  identification of a common Inuit founder mutation. Can. Med. Assoc. J. *187*, 102–107.
730  32. Maier, V., Jolicoeur, C., Rayburn, H., Takegahara, N., Kumanogoh, A., Kikutani, H.,
731  Tessier-Lavigne, M., Wurst, W., and Friedel, R.H. (2011). Semaphorin 4C and 4G are
732  ligands of Plexin-B2 required in cerebellar development. Mol. Cell. Neurosci. *46*, 419–431.
733  33. Kapur, S., Mehra, S., Gajjar, D., Vasavada, A., Kapoor, M., Sharad, S., Alapure, B.,
734  and Rajkumar, S. (2009). Analysis of single nucleotide polymorphisms of CRYGA and
735  CRYGB genes in control population of western Indian origin. Indian J. Ophthalmol. *57*,
736  197–201.
737  34. Perez-Oliva, A.B., Lachaud, C., Szyniarowski, P., Muñoz, I., Macartney, T., Hickson, I.,
738  Rouse, J., and Alessi, D.R. (2015). USP 45 deubiquitylase controls ERCC 1 – XPF
739  endonuclease-mediated DNA damage responses. EMBO J. *34*, 326–343.
740  35. Harris, K., and Nielsen, R. (2016). The genetic cost of neanderthal introgression.
741  Genetics Early online April 2, 2016; DOI: 10.1534/genetics.
742  36. Jorgensen, M.E., Bjeregaard, P., Borch-Johnsen, K., Backer, V., Becker, U.,
743  Jorgensen, T., and Mulvad, G. (2002). Diabetes and impaired glucose tolerance among
744  the Inuit population of greenland. Diabetes Care *25*, 1766–1771.
745  37. McNair, A., Gudmand-Hoyer, E., Jarnum, S., and Orrild, L. (1972). Sucrose
746  malabsorption in Greenland. Br. Med. J. *2*, 19–21.
747  38. Gudmand-Høyer, E., Fenger, H.J., Kern-Hansen, P., and Madsen, P.R. (1987).
748  Sucrase deficiency in Greenland. Incidence and genetic aspects. Scand. J. Gastroenterol.
749  *22*, 24–28.
750  39. Lim, E.T., Würtz, P., Havulinna, A.S., Palta, P., Tukiainen, T., Rehnström, K., Esko, T.,
751  Mägi, R., Inouye, M., Lappalainen, T., et al. (2014). Distribution and medical impact of
752  loss-of-function variants in the Finnish founder population. PLoS Genet. *10*, e1004494.
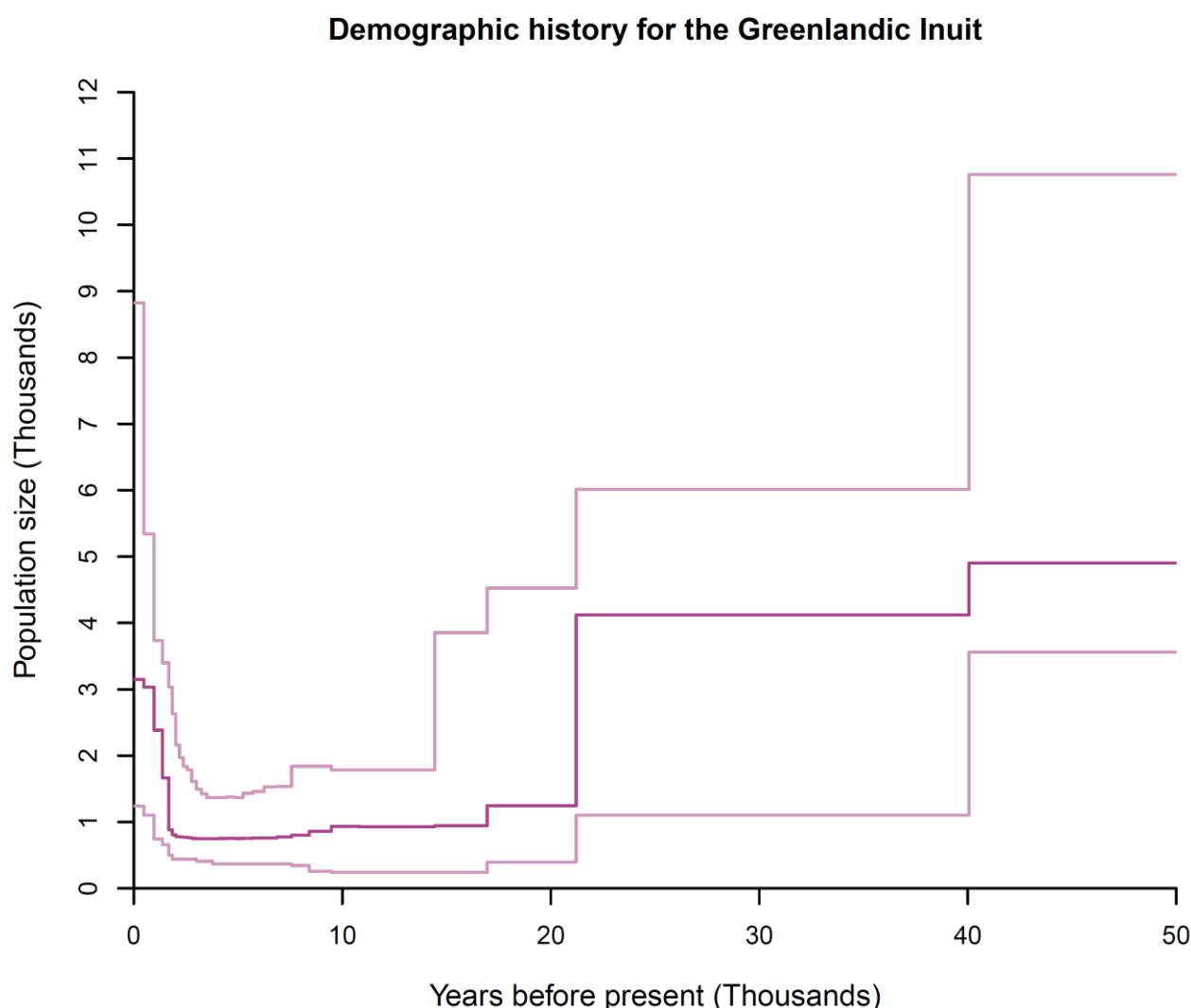753
754
755
756
757

## Figures



**Figure 1. Stairway plot for the Greenlandic Inuit population.**

The dark pink line shows the estimated diploid population size changes in discrete

increments for the last 50k years. The estimates were obtained with the method "Stairway

plot", which bases its estimates on the site frequency spectrum. The estimates are based

on an assumption of a mutation rate of $1.2 \times 10^{-8}$ per site per generation and a generation

time of 24 years. The light pink lines represent 95% CI based on bootstraps. This analysis

was based on 41.222.102 sites.

31

**Figure 2. Site frequency spectrum for six human populations.**

We used 11 randomly sampled individuals from each population to infer the site frequency

spectrum and excluded fixed categories. The GI population has fewer sites in the singleton

category, but more in the remaining more "common" categories. Each population is

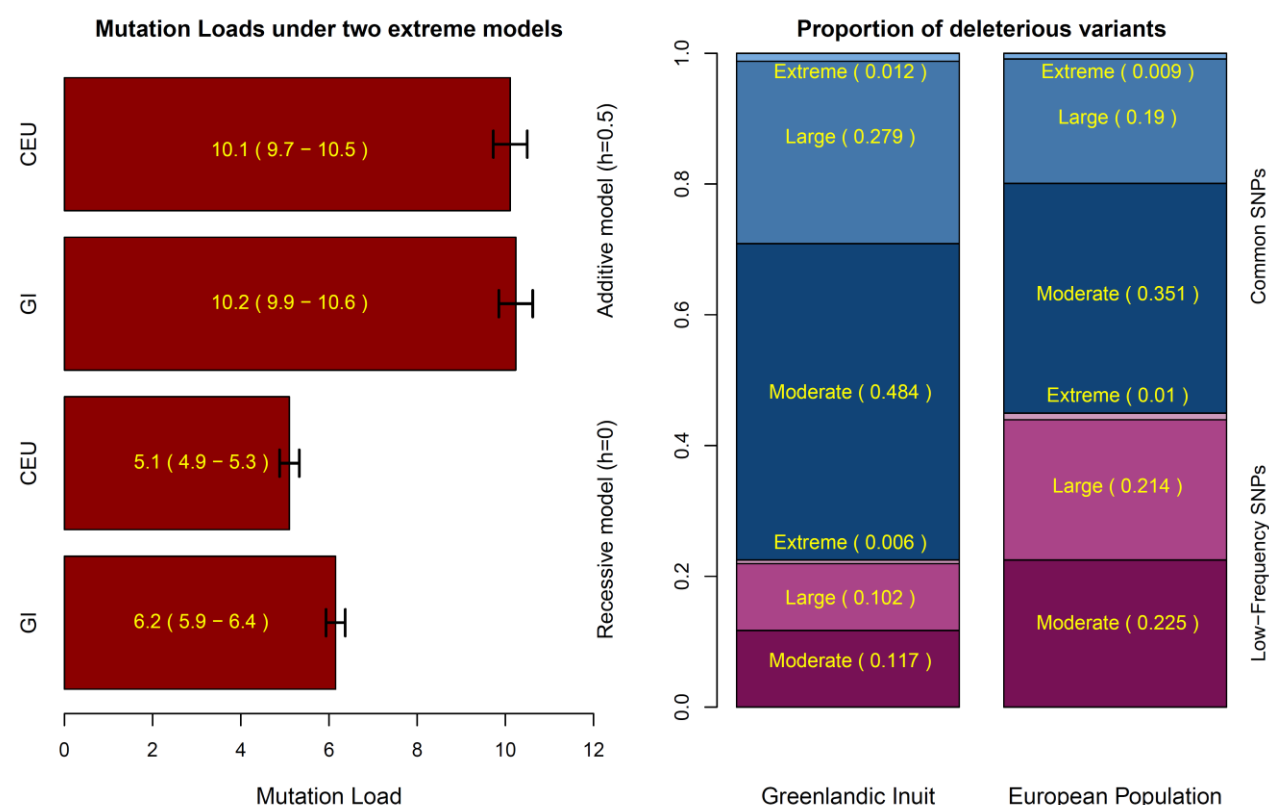followed by a $\pi_{var}$ estimate per variable site.

32

**Figure 3. Mutation Load and Load proportions.**

The left plot shows the genetic load using a fully additive model (top two bars) and a fully recessive model (bottom two bars). These genetic load values are based on annotated GERP scores converted to selection coefficients using the approach from Boyko et al.[6]. We note that if the selection coefficients from Henn et al.[5] are used instead, we see qualitatively similar results. The genetic load was calculated as in Lohmuller[11]. Black error bars indicate 95% C.I. The right plot shows the proportion of deleterious variants classified by GERP score (Moderate: 2<GERP<4, Large: 4<GERP<6, Extreme: 6<GERP) and their frequency in the populations (Low vs common). Low frequency SNPs are here defined as singletons and doubletons (equivalent of a frequency of at most 1/18≈0.056), while common SNPs are defined as tripletons or more than that including fixed derived sites (equivalent of a frequency of 1/12 or above).

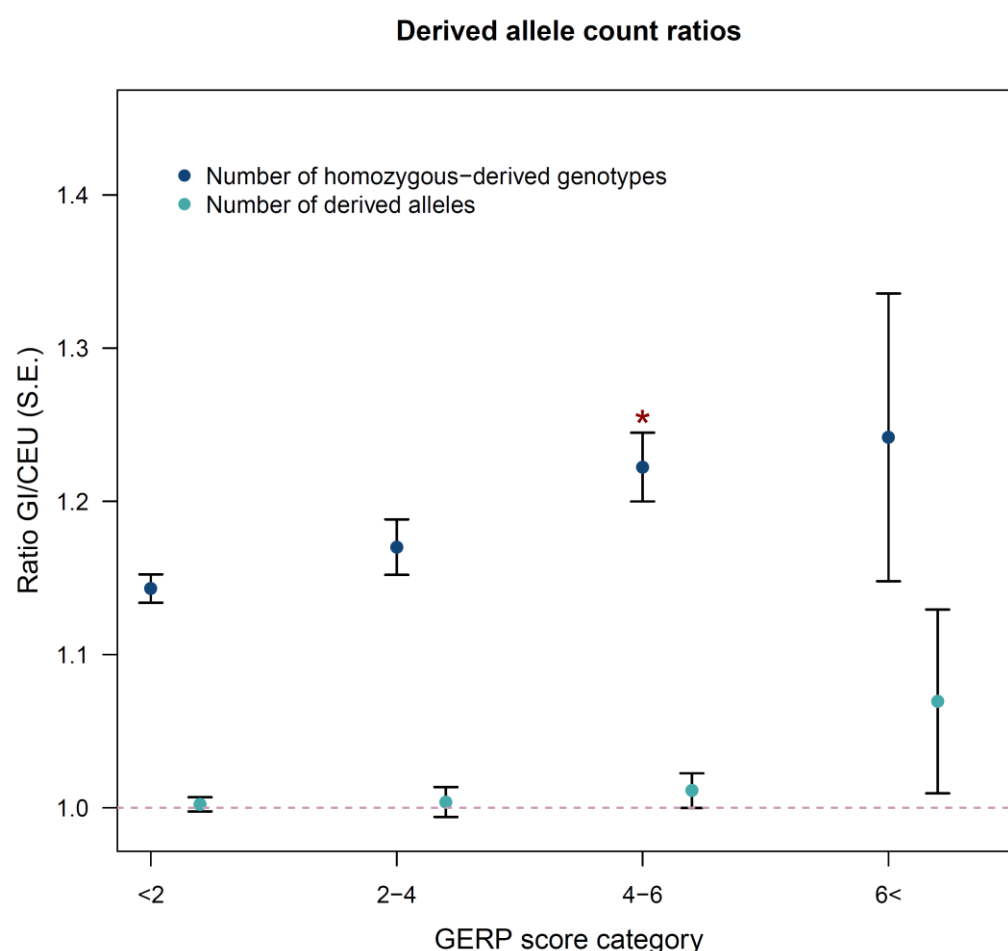**Derived allele count ratios**



792

**Figure 4. Load statistic ratios**.

For each of four GERP score categories two ratios calculated based on sites located in exons are shown: the ratio of derived allele counts in GI versus CEU (turquoise) and the ratio of homozygous-derived genotype counts in GI versus CEU (blue). The former can be viewed as an approximation for the ratio of load between the two populations under an additive model and latter can be viewed as an approximation for the ratio of load between the two populations under a recessive model. Standard error for each ratio is indicated by error bars. Additional Information is available in Table 1. The * indicates significance compared to the neutral GERP score category (GERP<2) ($P = 5.3 \times 10^{-4}$) using a $z$-test (compare column two for row 7 and 9 in Table 1).
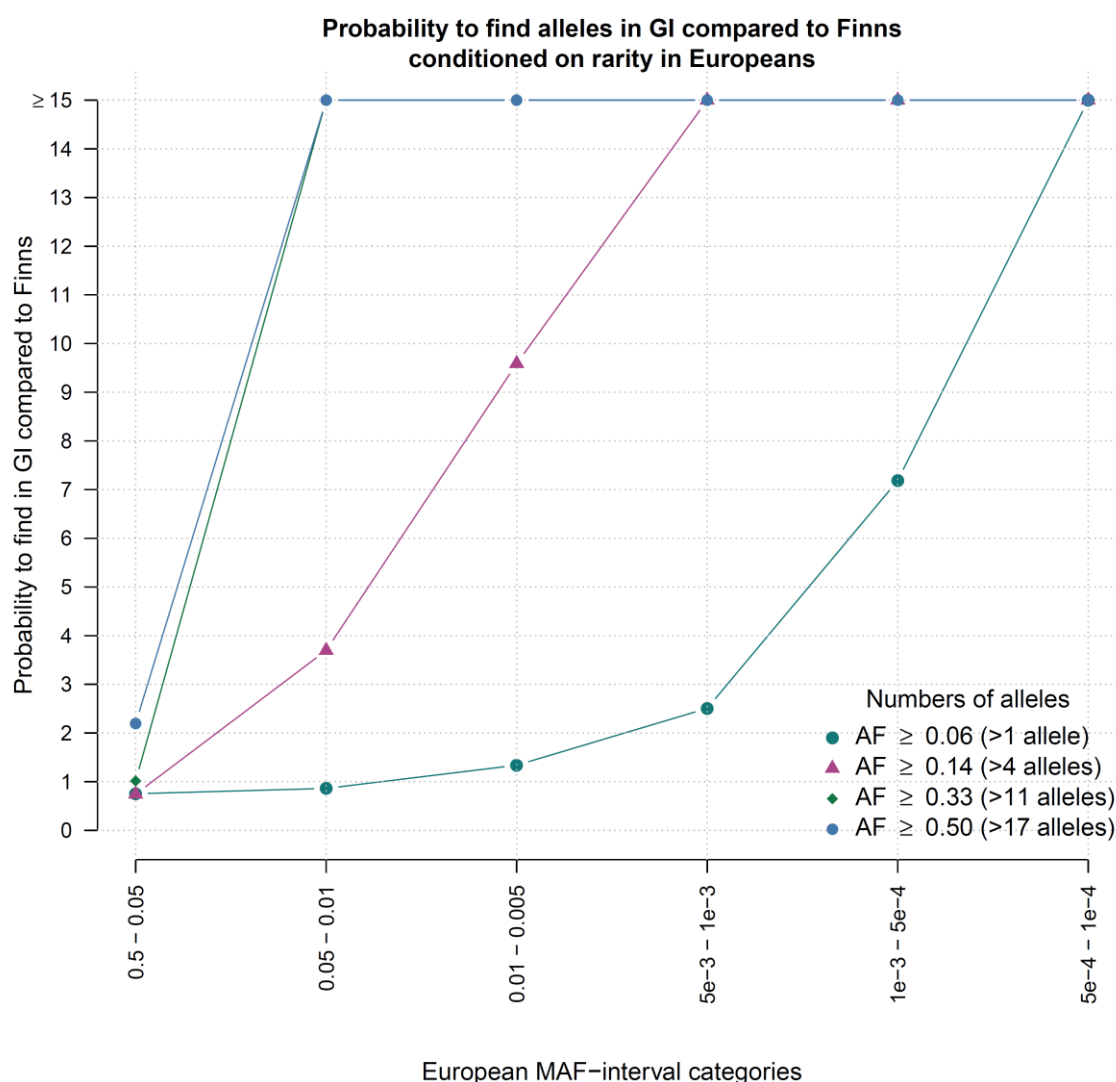
**Figure 5. Probability of common alleles.**

For alleles in different frequency categories in Europeans (x-axis) the points show how

much more likely the alleles are to be common in GI than in Finns (y-axis). Common is

here defined in 4 different ways, each represented by a specific color point, e.g. the dark

blue circles represents results for analyses made with common defined as more than 2 out

of 36 alleles, which corresponds to a frequency of ≈0.056 (for the remaining definitions,

see the figure legend).

# Tables

**Table 1. Summary of load investigations.**

Comparisons between 18 Greenlandic Inuit (GI) and 18 Utah Residents (CEPH) with Northern and

Western Ancestry (CEU) individuals using 93.047 SNPs sites located in genic regions with full

information for both populations. Ratios are the sum within GI divided by the sum within CEU.

Standard errors are given in parentheses.

| Category (N SNV) | Ratio of derived allele counts (SE) | Ratio of homozygous-derived genotype counts (SE) | Ratio of derived alleles in fixed sites (SE) | Derived alleles from fixed sites | Ratio of derived alleles per SNV (SE) | $R_{xy}$ (SE) |
|---|---|---|---|---|---|---|
| Derived Alleles (49631) | 1.002 (0.004) | 1.156 (0.008) | 4.413 (0.301) | GI:17.0% CEU:3.9% | 1.484 (0.009) | 1.004 (0.008) |
| Synonymous (20693) | 1.006 (0.005) | 1.156 (0.010) | 4.733 (0.420) | GI:18.2% CEU:3.9% | 1.420 (0.011) | 1.013 (0.012) |
| NCE+UTR (4919) | 1.006 (0.012) | 1.168 (0.022) | 5.268 (0.900) | GI:17.9% CEU:3.4% | 1.466 (0.023) | 1.013 (0.026) |
| Missense (23582) | 0.996 (0.006) | 1.151 (0.012) | 3.899 (0.377) | GI:15.5% CEU:4.0% | 1.550 (0.013) | 0.993 (0.013) |
| LoF (400) | 1.009 (0.063) | 1.296 (0.131) | 2.500 (2.840) | GI:9.3% CEU:3.8% | 1.625 (0.115) | 0.991 (0.104) |
| LoF[a] (278) | 1.045 (0.099) | 1.457 (0.286) | - | GI:0% CEU:4.4% | 2.026 (0.194) | 1.064 (0.142) |
| GERP < 2 (29183) | 0.999 (0.005) | 1.140 (0.009) | 4.293 (0.291) | GI:18.4% CEU:4.0% | 1.414 (0.009) | 0.999 (0.010) |
| GERP >2:<4 (9853) | 1.001 (0.010) | 1.169 (0.018) | 4.725 (0.757) | GI:15.3% CEU:3.1% | 1.517 (0.019) | 1.002 (0.019) |
| GERP >4:<6 (10211) | 1.011 (0.011) | 1.222 (0.022) | 4.659 (0.919) | GI:12.7% CEU:2.8% | 1.697 (0.023) | 1.021 (0.021) |
| GERP >6 (353) | 1.069 (0.060) | 1.244 (0.093) | - | GI:8.6% CEU:0% | 1.888 (0.130) | 1.127 (0.113) |

[a] = unambiguous sites (sites where all possible functional annotations were LOF annotations)

- indicates that too few sites fall in this the category to allow a meaningful result

836 **Table 2. Summary of Loss of Function mutations.**

837 LoF variant (Stop & splice variants) annotation for alleles that are rare in Europeans and East

838 Asians (below a MAF of 0.5%) but not lost.

| rs number | Gene name | Possible clinical implications | Genotypes among the 18 Inuit (Parents only) | Reference allele |
|---|---|---|---|---|
| rs12471298 | *SEMA4C* | Neonatal lethality | CC=12,CA=6 | C |
| rs61736969 | *TBC1D4* | Type 2 diabetes | GG=7,GA=7,AA=4 | G |
| rs116344874 | *CRYGA* | Cataract | GG=13,GA=5 | G |
| rs189281869 | *USP45* | DNA repair (following UV irradiation) | TT=15,TA=3 | T |
| rs189664399 | *SIGLEC7* | unknown | CC=15,CT=3 | C |
| rs201985967 | *SPAG4* | unknown | TT=13,TA=5 | T |

839
840

841

842

843

37