# 1 Limitations of lymphoblastoid cell lines

# 2 for establishing genetic reference

# 3 datasets in the immunoglobulin loci

4

5 Limitations of LCL for IG reference datasets

6

7 Oscar L. Rodriguez[1], Andrew J. Sharp[2], Corey T. Watson[1]

8

9 [1] Department of Biochemistry and Molecular Genetics, University of Louisville School of

10 Medicine, Louisville, KY USA

11 [2]Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New

12 York, NY, USA

13

14 Corresponding author: corey.watson@louisville.edu

15

# Abstract

Lymphoblastoid cell lines (LCLs) have been critical to establishing genetic resources for biomedical science. They have been used extensively to study human genetic diversity, genome function, and inform the development of tools and methodologies for augmenting disease genetics research. While the validity of variant callsets from LCLs has been demonstrated for most of the genome, previous work has shown that DNA extracted from LCLs is modified by V(D)J recombination within the immunoglobulin (IG) loci, regions that harbor antibody genes critical to immune system function. However, the impacts of V(D)J on data generated from LCLs has not been extensively investigated. In this study, we used LCL-derived short read sequencing data from the 1000 Genomes Project (n=2,504) to identify signatures of V(D)J recombination. Our analyses revealed sample-level impacts of V(D)J recombination that varied depending on the degree of inferred monoclonality. We showed that V(D)J associated somatic deletions impacted genotyping accuracy, leading to adulterated population-level estimates of allele frequency and linkage disequilibrium. These findings illuminate limitations of using LCLs for building genetic resources in the IG loci, with implications for interpreting previous disease association studies in these regions.

# Author summary

Lymphoblastoid cell lines (LCLs) are cells that have been manipulated to proliferate indefinitely in order to provide a replenishable source of DNA. However, because these cell lines are derived from B cells which have undergone V(D)J recombination they contain somatic deletions within regions of the genome that encode antibody genes. Although several large collaborative projects have utilized DNA from LCLs to generate invaluable genomic resources for the scientific community, the negative impacts of cell line artifacts in these regions of the genome

39    have not been fully appreciated. In this study, we used newly released sequencing data from a

40    large collection of LCLs to determine that the non-inherited artificial deletions within the antibody

41    gene loci can have detrimental effects on downstream genetic analyses.

42

43    **Keywords**: Immunoglobulin heavy chain locus; Lymphoblastoid cell lines; Genotyping; Genetic

44    Association Analysis; Linkage Disequilibrium

# Introduction

46         Lymphoid blastoid cell lines (LCL) are generated by infecting B cells with the Epstein

47    Barr Virus (EBV)[1] to create immortalized cell lines. Various consortia, including The

48    International HapMap Project[2, 3], 1000 Human Genome Project (1KGP)[4–6], Genome In A

49    Bottle[7, 8] and Human Genome Structural Variation Consortium[9]  have used DNA from LCLs

50    to characterize common genetic variation, generate gold standard sets of small insertions and

51    deletions (indels), and comprehensively genotype structural variants (SV). Variant call sets from

52    these initiatives have been instrumental to the genomics community, and are routinely used in

53    genome-wide association studies (GWAS) and other genetic studies. Genome-wide genotypes

54    from LCLs have been shown to be nearly identical to genotypes derived from whole blood or

55    peripheral blood mononuclear cells (PBMC) using SNP arrays[10], whole exome

56    sequencing[11, 12] and whole genome sequencing[13]. However, somatic LCL-associated

57    alterations are present in particular regions of the genome, namely within the immunoglobulin

58    (IG) heavy (IGH) and light (lambda, IGL; kappa, IGK) chain loci. These alterations could impact

59    sequencing, mapping, and genotype results in these regions, with potential implications for

60    downstream uses of these data.

61     The IG loci encode the variable (V), diversity (D) and joining (J) gene segments that

62     serve as the building blocks for the expression of functional B cell receptors (BCRs) and

63     antibodies (Abs). During B cell development, the V, D, and J gene segments within each IG

64     locus (V and J in the case of IGL and IGK) are somatically rearranged through a process called

65     V(D)J recombination[14]. During this process, intervening DNA between recombined V, D, and J

66     segments is excised. The size of these somatic deletions on the recombined chromosome

67     depends on the selected V, D, and J genes, but can extend 100's of Kb, and will vary from cell

68     to cell. Collectively, DNA isolated across a pool of B cells (e.g., naive B cells) representing many

69     independent V(D)J recombination events would be expected to represent each germline

70     haplotype present in a given sample (Fig. 1). In contrast, a pool of B cells originating from a

71     single or dominant expanded B cell would harbor DNA not fully representative of both paternal

72     and maternal germline haplotypes within the IG loci (Fig. 1). In the latter instance, genotyping

73     methods dependent on different read alignment signatures such as read depth/coverage,

74     discordant read mapping, soft-clipped or split reads could produce inaccurate germline

75     genotypes.

76     Recent long read sequencing and assembly of complete IGH haplotypes from selected

77     1KGP individuals has revealed the presence of V(D)J recombination associated deletions[15],

78     indicating that genotypes derived from such samples within regions impacted by V(D)J

79     recombination are inaccurate. While it has been speculated previously that V(D)J recombination

80     would have negative impacts on LCL-derived sequencing data [16–18], this has not been

81     comprehensively investigated. Given this, we sought to evaluate the extent of sample-level

82     V(D)J recombination in LCL-derived short read sequencing data from the 1KGP, and assess

83     downstream impacts of these somatic events. We demonstrate that short read data is affected

84     by V(D)J recombination and, depending on the sample, is derived from either single dominant or

85     multiple B cell clones. We show that variation in sample clonality is associated with variability in

86    genotyping accuracy, negatively impacting estimates of allele frequency and linkage

87    disequilibrium (LD). These data raise important considerations for using 1KGP genotypes to

88    augment genetic association studies in the IG loci, and in addition to other issues discussed

89    previously[16–18], may further explain the paucity of disease associations within these complex

90    regions of the genome.

# Materials/Subjects and Methods

## 1000 Human Genome Project Data

93    Paired-end 150 bp PCR-free 30X coverage Illumina data on 2504 individuals from the 1KGP[19]

94    was downloaded from the European Bioinformatics Institute (EBI) under the study ID

95    ERP114329. 1KGP phase 3[4] SNPs were downloaded from

96    ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20190

97    312_biallelic_SNV_and_INDEL/ALL.chr14.shapeit2_integrated_snvindels_v2a_27022019.GRC

98    h38.phased.vcf.gz.

## Assessment of insert sizes, read depth, and the identification of

## V(D)J recombination events

101    Insert sizes for read pairs within each sample were calculated using the Picard

102    (https://github.com/broadinstitute/picard) CollectInsertSizeMetrics tool with the following

103    parameters: `--DEVIATIONS 1000000 --MINIMUM_PCT 0` for reads spanning

104    chr14:105862198-107043718. The same tool and parameters were used for ten random 1.2 MB

105    windows across the genome. To calculate read depth across the IGHD and V regions, we used

106    samtools[23] (IGHD, hg38, chr14:105,865,458-105,939,756; IGHV, hg38, chr14:105,939,756-

107   106,883,718). To identify read pairs within each sample representing V(D)J recombination

108   events, we utilized a custom python script. For each sample, the number of clones was

109   calculated by counting the number of unique IGHV and IGHJ gene pairs detected. The

110   frequency of each "clone" was calculated by determining the number of reads mapping to a

111   unique IGHV and IGHJ gene pair, and taking this as a fraction of the total number of reads

112   assigned to any IGHV/IGHJ pair. The sizes of somatic deletions were determined by calculating

113   the genomic distances between the IGHJ and IGHV genes utilized in a given V(D)J event.

## Analysis of heterozygosity, allele frequency, and LD

115   The percentage of heterozygous SNPs was calculated for each sample in the

116   centromeric and telomeric region of the selected V gene in the dominant clone. Two VCFs were

117   created for each region using tabix to select the region of interest on the 1KGP phase 3 VCF.

118   The telomeric region was always set to start at the beginning of *IGHV6-1* (chr14:105,939,756;

119   GRCh38). The fraction of heterozygous SNPs was calculated by counting the number of

120   heterozygous SNPs over the total number of SNPs in each VCF.

121   To assess the effects of V(D)J recombination on allele frequency estimates and LD, we

122   subsetted the 1KGP phase 3 genotype call set to samples from the "AFR" superpopulation,

123   further selecting samples representing two extremes of clonality (0-25%, n=38; 75-100%, n=38).

124   Allele frequencies for each set of samples was calculated using the vcftools `freq` tool. The LD

125   scores were calculated for the African samples of each set of samples using the vcftools `hap-

126   r2` tool with the parameters `--ld-window 1000000 --min-r2 0.01`.

127

# Results

## Detecting signatures of V(D)J recombination

130     To determine the effect of V(D)J recombination on whole genome sequencing (WGS)

131  data in IGH, we used paired-end 150 bp PCR-free Illumina data on 2,504 individuals from the

132  1KGP, recently resequenced to high coverage[19]. The occurrence of V(D)J recombination

133  results in large somatic deletions within the IGH locus spanning the IGHJ, IGHD, and IGHV

134  regions. To assess signatures of these somatic deletions we first analyzed paired-end mapping

135  distances, as measured by the predicted "insert sizes". We reasoned that the presence of V(D)J

136  recombination events would result in larger insert sizes, and that these would be enriched within

137  IGH. To assess this, we calculated the number of read pairs with an insert size >900 bp (two

138  times the library DNA insert size) at 10 random 1.2 MB windows (the length of the IGH locus in

139  GRCh38) across the genome from five individuals chosen at random. Across these regions in

140  the selected individuals, we observed that 0.08% to 0.12% (mean = 0.10%) of the paired-end

141  reads contained an insert size greater than 900 bps. In contrast, across all samples, 0.13% to

142  1.55% (mean = 0.49%) of paired-end reads in IGH contained an insert size greater than 900

143  bps. This is almost a 5-fold increase in the number of paired-end reads with a larger insert size

144  (Fig. 2A).

145     To further evaluate the effect of V(D)J recombination, we calculated the coverage over

146  the IGHD region. During B cell development, through the formation of the pre-B cell receptor,

147  V(D)J recombination results in a loss of DNA between selected IGHJ and IGHD gene segments

148  on each homologous chromosome[20]. Therefore, if V(D)J recombination has occurred, there

149  should be limited to no coverage within the IGHD region. The closest IGHJ and IGHD gene pair,

150  *IGHJ1* and *IGHD7-27*, are within 100 bps, so we assessed coverage between *IGHJ1* and

151    *IGHD1-26* (~15 Kb) across all 2504 individuals, and observed a mean coverage of 0.76X (range

152    = 0 - 30.6). We also observed a mean coverage of 7.30X between *IGHJ1* and *IGHV6-1*, in

153    contrast to a mean coverage of 24.68X across the entirety of the IGHV region (Fig. 2B).

154        We next used the paired-end data to directly detect V(D)J recombination events by

155    identifying read pairs with one mate overlapping an IGHJ gene and the other mate overlapping

156    an IGHV gene. We detected read pairs representing V(D)J recombinants in all 2,504 samples,

157    with a mean of 22 V(D)J associated read pairs per sample (range = 2 - 66). V(D)J

158    rearrangements utilizing *IGHJ4* and *IGHV3-23* were the most common (Supplemental Fig. 1).

159    Taken together, these three pieces of evidence, increased insert sizes, decreased coverage

160    over IGHD, and the direct detection of read pairs overlapping V(D)J recombination events,

161    indicate that in fact LCLs utilized by the 1KGP cohort have undergone V(D)J recombination.

## Sequencing data derived from multiple B cell clones

163    Given that V(D)J recombination has occurred across the cohort, we sought to determine

164    whether all samples were affected equally. We reasoned that samples with sequencing data

165    from a single B cell clone (monoclonal) or from multiple clones (polyclonal) will be differentially

166    impacted by the effects of V(D)J (Fig. 3A). We therefore sought to determine the number and

167    frequency of V(D)Js in each sample. To do this, we assigned read pairs overlapping V(D)J

168    events to their respective combination of IGHJ and IGHV genes. Reads across a given dataset

169    harboring the same IGHJ/IGHV combination were grouped, and used as a proxy for a group of

170    clonally related sequences. We thus took the number of grouped sequence reads to represent

171    the frequency of a particular IGHJ/IGHV combination, which we heretofore refer to as a clone

172    (Fig. 3A). Following this, we calculated  the number of unique IGHJ/IGHV combinations

173    ("clones") present in each sample, and their relative frequency, allowing us to approximate the

174    number of different B cell clones represented in a sample. We found that sequences across

175    samples in the cohort were derived from a mean of 10.76 B cell clones (range = 1 - 30; Fig. 3B).

176    From this, 18 samples were predicted to be monoclonal, represented by reads mapping to only

177    a single clone. We also reasoned that polyclonal samples, represented by many clones, but in

178    which the majority of sequencing data is predicted to be derived from a dominant clone will have

179    profiles similar to those observed for monoclonal samples. To estimate this, we asked what

180    proportion of all sequences containing a V(D)J recombination event were represented by the

181    most frequently observed clone. In doing this, we found that in 407 and 88 samples,

182    respectively, 50% and 75% of all reads containing V(D)J recombination events mapped to a

183    single clone. This indicated that although these samples had a polyclonal signature, the majority

184    of sequencing data was likely derived from a single dominant clone. Based on this

185    approximation, across samples, the top clone identified contributed on average 32.94% (range

186    = 5 - 100%) of the sequencing data (Fig. 3C). We observed a modest population bias, with

187    African and European individuals containing a slightly greater fraction of sequencing data from a

188    single clone (Supplemental Fig. 2).

189         Depending on the clonality of the sample and the genes involved in the primary V(D)J

190    recombination event present within a sample, the size of the region impacted was expected to

191    vary. We estimated this in each sample based on the most prevalent IGHJ/IGHV gene

192    combination observed, revealing that in many samples the predicted size of these somatic

193    deletions was extensive. In the 18 monoclonal samples, we found that an average of 464 Kb of

194    the IGH locus was impacted by V(D)J recombination (range = 74.3 - 937.2 Kb). In samples that

195    were polyclonal, but still represented by a dominant clone (i.e., those in which >50% and >75%

196    of sequence data were derived from a single IGHJ/IGHV combination), the regions impacted by

197    V(D)J recombination were found to be on average 356 Kb (50%, range = 74.3 - 945.0 Kb) and

198    401 Kb (75%, range = 74.3 - 945.0 Kb) in size (Fig. 3D). Regions affected by V(D)J

199    recombination in each sample are provided in Supplementary Table 1.

# The effects of V(D)J recombination on genotype call sets

200       We reasoned that V(D)J events could impact the accuracy of sample- and population-

201

202    level genotypes in two primary ways: 1) the loss of DNA and reduced read coverage over

203    extended regions of the locus would result in the increased likelihood of calling homozygous

204    genotypes at heterozygous positions; 2) somatic hypermutations (SHMs) in recombined genes

205    could introduce false heterozygous SNPs. Furthermore, these effects would likely be more

206    prominent in monoclonal samples, as well as polyclonal samples represented by dominant

207    clones.

208       To investigate these potential impacts, we first evaluated the number of heterozygous

209    SNPs in the centromeric and telomeric regions of the recombined IGHV gene in monoclonal

210    samples. We found that the mean percentage of heterozygous variants telomeric of the IGHV

211    gene used for V(D)J recombination was 3.4 fold higher than the mean percentage of

212    heterozygous variants centromeric of the IGHV gene ($P = 0.003$, two-sided paired Wilcoxon

213    test; Fig. 4A). To assess this effect in polyclonal samples, we split individuals into four groups

214    representing varying degrees of clonal bias, based on whether 0-25%, 25%-50%, 50%-75% and

215    75%-100% of sequencing data within a given sample came from the dominant clone. In

216    samples with 0 to 25% of sequencing data from the dominant clone, for which we expected to

217    observe minimal impacts on genotyping, the mean percentage of heterozygous variants

218    telomeric (62%) of the selected recombined IGHV gene was 1.04-fold higher than in the

219    centromeric region (59%; $P$=0.006, two-sided paired Wilcoxon test; Fig. 4B). We noted

220    significant differences in the remaining three groups as well (Fig. 4B), but the average fold-

221    differences between heterozygous percentages telomeric to the V(D)J event relative to

222    centromeric to the V(D)J event were greater, and was greatest in samples from the 75%-100%

223    group (2.08-fold).

224        Additionally, we evaluated the number of heterozygous SNPs overlapping IGHV genes

225    most frequently selected for V(D)J recombination in each sample. When stimulated by an

226    antigen, B cells acquire SHMs within IG V, D, and J genes as a means to increase antibody

227    affinity[21]. Therefore, SHMs in the recombined IGHV gene are more likely to be detected in

228    monoclonal or polyclonal samples with reads primarily derived from a dominant clone. Indeed,

229    there was a significant positive correlation (R = 0.33, p-value < 2.2e-16) between the

230    contribution of sequencing data from the dominant clone and the number of the heterozygous

231    variants within the IGHV gene selected by V(D)J recombination (Supplementary Fig. 3A). We

232    also directly compared the number of heterozygous genotypes within the recombined IGHV

233    genes to non-recombined IGHV genes across all samples, and observed an average of 1.92

234    heterozygous positions in the recombined IGHV genes, compared to 0.5 in the non-recombined

235    IGHV genes (Supplementary Fig. 3B).

# 236 The effects of V(D)J recombination on estimates of allele

## 237 frequency and linkage disequilibrium

238        The previous section detailed the effects on genotypes due to V(D)J recombination.

239    Given that genotypes are used to determine allele frequencies in a population, we set out to test

240    if allele frequencies differed between samples that are more or less monoclonal. The allele

241    frequencies of common SNPs (MAF > 0.05) were compared between samples within

242    superpopulations with 0-25% (less monoclonal) and 75-100% (more monoclonal) of sequencing

243    data derived from the dominant clone. Of the 4,354 SNPs analyzed, 1,258 (29%) had an allele

244    frequency difference greater than 0.05 (Fig. 5A). Since V(D)J recombination excises DNA 3' of

245    selected IGHV genes, we would expect to observe more genotyping errors caused by V(D)J

246    related somatic deletions within the centromeric region of the IGHV locus. Consistent with this,

247    we observed greater differences in allele frequencies within the proximal (centromeric) region of

248    the locus when comparing estimates generated from less monoclonal samples to more

249    monoclonal samples (Fig. 5B).

250        Genotypes are also used for the calculation of LD between SNPs. Given the

251    demonstrated impact on allele frequency estimates, we reasoned that effects on genotype

252    accuracy would also impact LD estimates. To assess this, we chose 76 samples from the

253    African superpopulation representing extremes of clonality. Two groups of 38 samples each

254    from the lower clonality group and the higher clonality group were selected. The LD $r^2$ values

255    across the locus were computed and compared between the two groups (Fig. 5C), revealing

256    different LD structure. We found that 11% (236,827) of the SNP pairs exhibited differences in

257    LD ($r^2$) greater than 0.1 (Supplementary Fig. 4). The differences in allele frequencies and LD

258    estimates observed here indicated that inaccurate genotypes resulting from impacts of V(D)J

259    recombination also affect downstream analyses.

# Discussion

261        Previous studies have concluded that there are minimal differences between genotypes

262    from matched LCL and non-LCL samples[10–13]. While true on a genome-wide scale, here we

263    show that the impact on the IGH locus is more apparent due to V(D)J recombination. Using

264    1KGP samples (n=2504) recently resequenced on the Illumina NovaSeq platform to 30x

265    coverage using PCR-free 2x150 bp libraries, we evaluated different sequencing features

266    affected by V(D)J recombination and SHM within the IGH locus. Specifically, we demonstrated

267    that signatures of V(D)J recombination within LCL-derived DNA can be observed, including

268    increased insert sizes of read mate-pairs, decreased read coverage over the IGHD and

269    proximal IGHV gene regions, as well as direct evidence of somatically recombined IGHJ and

270    IGHV genes. By assessing the frequency of specific IGHJ/IGHV recombination events within

271    each sample, we were able to estimate the number of approximate B cell clones likely

272    represented within a sample, and determine the proportion of sequencing data derived from

273    each B cell clone, revealing variation in clonality across samples. Importantly, we were able to

274    determine that V(D)J recombination can result in loss of DNA spanning large segments of the

275    locus, with clear impacts on variant genotyping. The extent of these effects varied between

276    samples based on the gene segments involved in the primary V(D)J recombination event, and

277    the degree of monoclonality observed. Together these observations highlight critical limitations

278    of using LCLs to develop comprehensive reference resources for the IGH locus at the sample

279    and population level.

280        It has previously been argued that the locus complexity of IGH has made it difficult to

281    study using high-throughput approaches such as short-read data and genotyping arrays[15–17].

282    This has impeded our ability to accurately characterize genetic diversity within IGH, and robustly

283    test hypotheses about the functional role of IGH germline variation in disease risk and antibody-

284    mediated immunity. The analyses we have conducted here indicate that the large-scale use of

285    LCLs for establishing genetic reference panels in IGH may also present additional barriers to

286    effectively interrogating IGH in genetic studies with downstream implications that need to be

287    considered. For example, LCL-derived datasets such as the 1KGP have been critical for

288    establishing population-genetic metrics across the genome, and have been used to augment

289    GWAS and inform functional and population genetic studies. For example, consortia efforts

290    such as gnomAD[22] have aggregated data from multiple sources, including LCL-derived data

291    from the 1KGP, to power such studies. However, we have shown here that genotype, allele

292    frequency, and LD estimates are incorrect for much of IGH due in part to impacts of V(D)J

293    events in the data. This highlights a need to reconsider use of these cohorts for such purposes.

294    We argue that, at a minimum, the continued use of LCL-derived datasets could be improved by

295    removing erroneous genotypes caused by V(D)J recombination induced deletions. As part of

296     this study, we have released a BED file with the coordinates of V(D)J recombined induced

297     deletions for each sample (Supplementary Table 1). It is possible that the development of

298     genotyping pipelines that account for such data anomalies on a per-sample basis would lead to

299     more accurate estimates of genotype and allele frequencies within IGH, with potential

300     downstream implications for improving imputation approaches utilized by GWAS. Finally, while

301     the focus of our study has been on the IGH locus, these observations would be applicable to the

302     IGL and IGK loci as well.

# 303 Acknowledgements

# 308 Competing Interests

309     The authors declare no competing interests.

# 310 References

311     1.  Frisan T, Levitsky V, Masucci M Generation of Lymphoblastoid Cell Lines (LCLs). Epstein
312         Barr Virus Protocols 125–127

313     2.  Consortium TIH, The International HapMap Consortium (2005) A haplotype map of the
314         human genome. Nature 437:1299–1320

315     3.  International HapMap Consortium, Frazer KA, Ballinger DG, et al (2007) A second
316         generation human haplotype map of over 3.1 million SNPs. Nature 449:851–861

317     4.  1000 Genomes Project Consortium, Auton A, Brooks LD, et al (2015) A global reference for
318         human genetic variation. Nature 526:68–74

319   5.  1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD, DePristo MA,
320        Durbin RM, Handsaker RE, Kang HM, Marth GT, McVean GA (2012) An integrated map of
321        genetic variation from 1,092 human genomes. Nature 491:56–65

322   6.  Sudmant PH, Rausch T, Gardner EJ, et al (2015) An integrated map of structural variation
323        in 2,504 human genomes. Nature 526:75–81

324   7.  Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, Salit M (2014)
325        Integrating human sequence data sets provides a resource of benchmark SNP and indel
326        genotype calls. Nat Biotechnol 32:246–251

327   8.  Zook JM, Catoe D, McDaniel J, et al (2016) Extensive sequencing of seven human
328        genomes to characterize benchmark reference materials. Sci Data 3:160025

329   9.  Chaisson MJP, Sanders AD, Zhao X, et al (2019) Multi-platform discovery of haplotype-
330        resolved structural variation in human genomes. Nat Commun 10:1784

331  10.  Herbeck JT, Gottlieb GS, Wong K, Detels R, Phair JP, Rinaldo CR, Jacobson LP, Margolick
332        JB, Mullins JI (2009) Fidelity of SNP array genotyping using Epstein Barr virus-transformed
333        B-lymphocyte cell lines: implications for genome-wide association studies. PLoS One
334        4:e6915

335  11.  Londin ER, Keller MA, D'Andrea MR, Delgrosso K, Ertel A, Surrey S, Fortina P (2011)
336        Whole-exome sequencing of DNA from peripheral blood mononuclear cells (PBMC) and
337        EBV-transformed lymphocytes from the same donor. BMC Genomics 12:464

338  12.  Schafer CM, Campbell NG, Cai G, et al (2013) Whole exome sequencing reveals minimal
339        differences between cell line and whole blood derived DNA. Genomics 102:270–277

340  13.  Nickles D, Madireddy L, Yang S, Khankhanian P, Lincoln S, Hauser SL, Oksenberg JR,
341        Baranzini SE (2012) In depth comparison of an individual's DNA and its lymphoblastoid cell
342        line using whole genome sequencing. BMC Genomics 13:477

343  14.  Jung D, Alt FW (2004) Unraveling V(D)J recombination. Cell 116:299–311

344  15.  Rodriguez OL, Gibson WS, Parks T, et al A novel framework for characterizing genomic
345        haplotype diversity in the human immunoglobulin heavy chain locus.
346        https://doi.org/10.1101/2020.04.19.049270

347  16.  Watson CT, Breden F (2012) The immunoglobulin heavy chain locus: genetic variation,
348        missing data, and implications for human disease. Genes Immun 13:363–373

349  17.  Watson CT, Matsen FA, Jackson KJL, Bashir A, Smith ML, Glanville J, Breden F, Kleinstein
350        SH, Collins AM, Busse CE (2017) Comment on "A Database of Human Immune Receptor
351        Alleles Recovered from Population Sequencing Data." The Journal of Immunology
352        198:3371–3373

353  18.  Watson CT, Steinberg KM, Huddleston J, et al (2013) Complete haplotype sequence of the
354        human immunoglobulin heavy-chain variable, diversity, and joining genes and
355        characterization of allelic and copy-number variation. Am J Hum Genet 92:530–546

356  19.  Byrska-Bishop M, Evani US, Zhao X, et al (2021) High coverage whole genome
357        sequencing of the expanded 1000 Genomes Project cohort including 602 trios. bioRxiv

358        2021.02.06.430068

359    20. Vettermann C, Schlissel MS (2010) Allelic exclusion of immunoglobulin genes: models and
360        mechanisms. Immunol Rev 237:22–42

361    21. Murphy KM, Weaver C (2016) Janeway's Immunobiology. Garland Science, Taylor &
362        Francis Group, LLC

363    22. Karczewski KJ, Francioli LC, Tiao G, et al (2020) The mutational constraint spectrum
364        quantified from variation in 141,456 humans. Nature 581:434–443

365    23. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin
366        R, 1000 Genome Project Data Processing Subgroup (2009) The Sequence Alignment/Map
367        format and SAMtools. Bioinformatics 25:2078–2079

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382 # Figures



383

384 **Figure 1. Generation of LCLs can lead to loss of sequencing data**

385 LCLs are generated from a pool of B cells. A pool of polyclonal B cells will contain different

386 V(D)J recombination events and therefore the collection of different B cell clones would lead to

387 no loss of sequencing data from either haplotype. LCLs generated from monoclonal B cells, *i.e.*

388 B cells with the same single V(D)J event, would lead to loss haplotype-specific sequencing data.

389

390

391

**Figure 2. Signatures of V(D)J recombination in paired-end WGS data**

(A) Percentage of paired-end read pairs with insert size greater than 900 bps in IGH and random genome-wide 1 Mb windows.

(B) WGS coverage between (1) *IGHJ1*, the IGHJ gene closest to the telomeric end and *IGHD1-26*, the second closest gene to the IGHJ gene cluster, (2) *IGHJ1* and *IGHV6-1*, the IGHV gene closest to the IGHD and IGHJ gene cluster and (3) IGHV region

407

**Figure 3. Clonality in different samples**

(A) IGV screenshots showing three samples with different clonality but have a dominant

clone with *IGHV4-39*. Red and gray lines represent large insert lengths. HG03397 has

25 read pairs aligning with to an IGHJ and IGHV gene, all of which align to *IGHV4-39*

and hence labelled as monoclonal. NA19131 and HG04180, which are polyclonal, have

25/33 and 14/26 read pairs aligning to *IGHV4-39*.

(B) Number of samples with different clonalities

(C) Percentage of sequencing data derived from dominant clone

(D) Amount of IGHV DNA lost in monoclonal samples and polyclonal samples with 75% and

50% of their sequencing data derived from the dominant clone

**Figure 4. Difference in heterozygosity in centromeric and telomeric regions of selected V(D)J recombination IGHV gene and clonality bias**

(A) Eighteen samples were predicted to be monoclonal. The percentage of heterozygous SNPs in the centromeric and telomeric regions of the IGHV gene selected for V(D)J recombination was calculated. Each line represents a single sample connecting the percentage of heterozygous SNPs in the centromeric and telomeric regions.
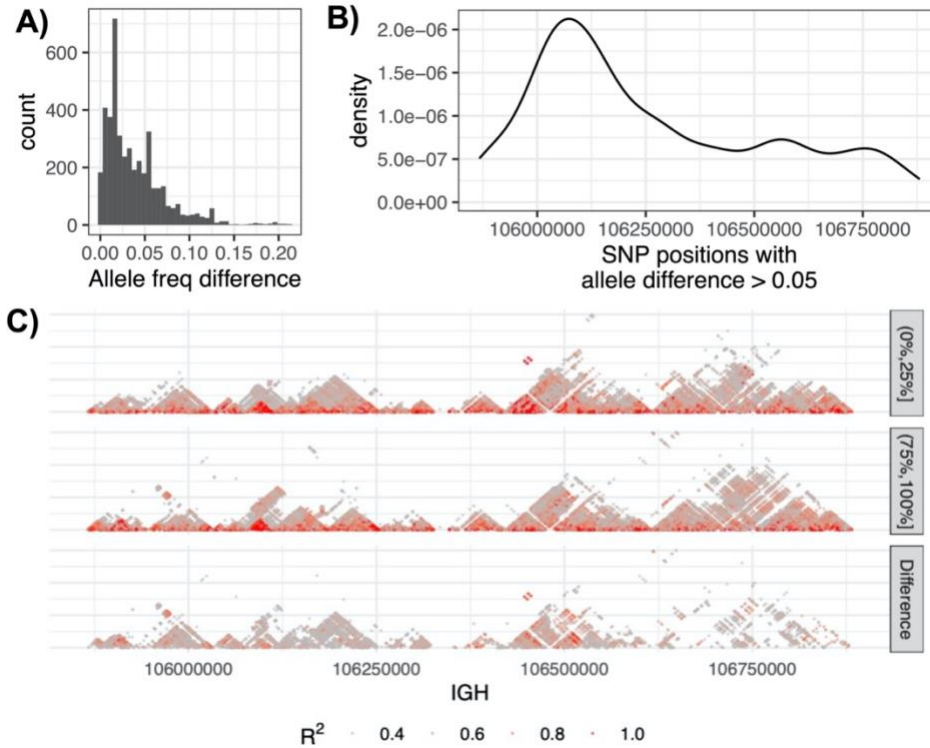
(B) Samples were split based on their clonality bias. Samples in the "(0,25]" group, the least clonal group, derived up to 25% of sequencing data from the dominant clone and samples in the "(75,100]" group, most clonal group, derived 75% to less than 100% of their sequencing data from a dominant clone. Similar to (A), the proportion of heterozygous SNPs was calculated in the centromeric and telomeric regions of the IGHV gene selected for V(D)J recombination.

436

**Figure 5. Allele frequency and LD differences in individuals with low and high clonality**

(A) The distribution of common allele frequency difference between individuals with low

clonality and high clonality, defined as 0 to 25% and 75% up to 100% of sequencing

data was derived from a single clone, respectively.

(B) Position in IGH with common allele frequency differences greater than 0.05

(C) LD for African individuals with low ("(0%,25%]") and high clonality ("(75%,100%]"), and
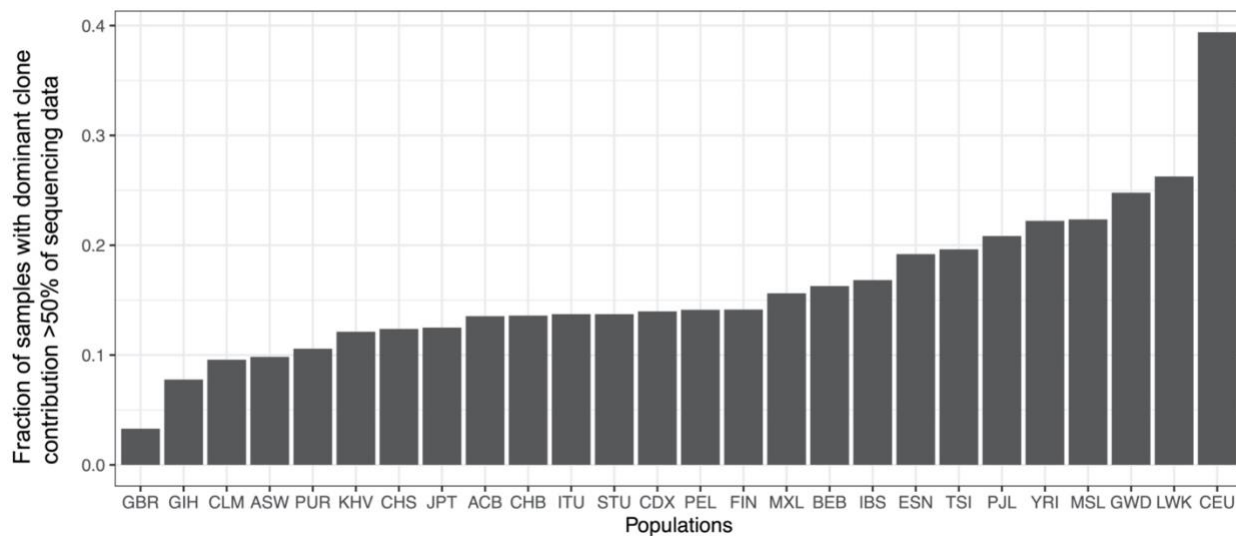
the difference in LD between both groups.

448    # Supplementary Figures



449

**Supp Fig 1. Number of samples with dominant clone containing IGHV and J pair**

451

**Supp Fig 2. Fraction of samples per population where the dominant clone is more than**

**50% and 75% prevalent**
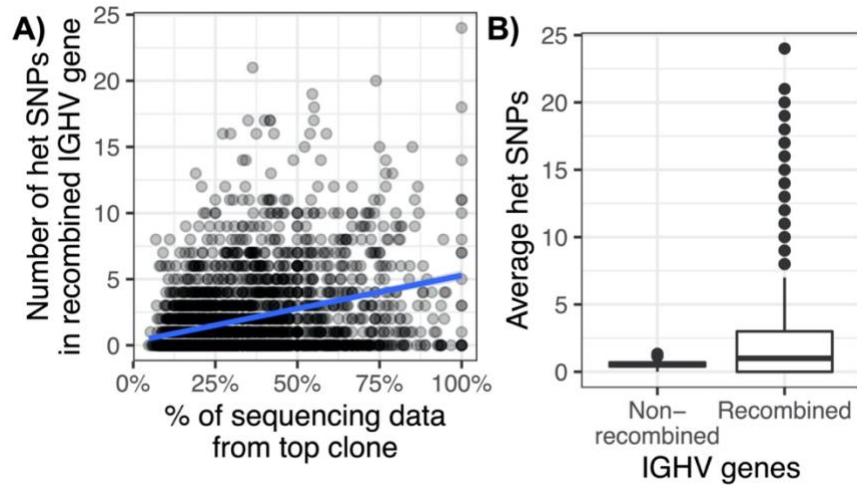
452

453

454

455

456

457

458

459

**Supp Fig 3. Effect of V(D)J recombination and clonality on heterozygous SNP calling**

460

461    (A) The number of heterozygous SNP in the V(D)J IGHV selected gene compared to the

462    percentage of sequencing data from dominant clone

463    (B) The average number of heterozygous SNPs between IGHV genes selected for V(D)J

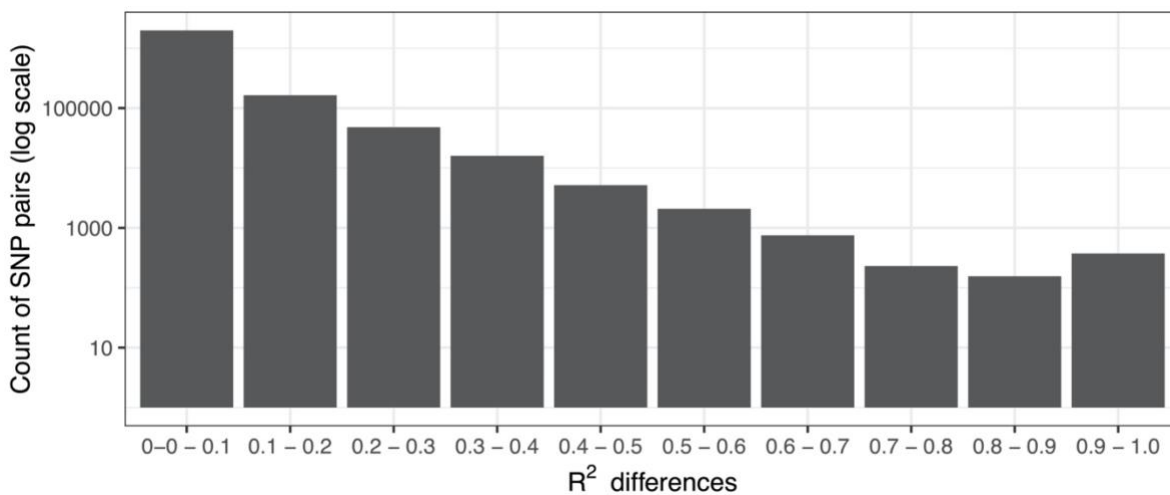464    recombination and not selected

465

466

467

468

469

470

471

472

473

474

**Supp Fig 4. LD difference between African individuals with low clonality (0-25%) and high clonality (75%-100%).**