# Polymorphisms in immunoglobulin heavy chain variable genes and their upstream regions

Ivana Mikocziova[1][†][*], Moriah Gidoni[2][†], Ida Lindeman[1], Ayelet Peres[2], Omri Snir[1], Gur Yaari[2][‡], Ludvig M. Sollid[1][‡]

[1] K.G.Jebsen Centre for Celiac Disease Research and Department of Immunology, University of Oslo and Oslo University Hospital, 0372 Oslo, Norway
[2] Faculty of Engineering, Bar Ilan University, Ramat Gan 5290002, Israel

† Joint First Authors

‡ Joint Last Authors

* To whom correspondence should be addressed. Email: ivana.mikocziova@medisin.uio.no

**ABSTRACT**

Germline variations in immunoglobulin genes influence the repertoire of B cell receptors and antibodies, and such polymorphisms may impact disease susceptibility. However, the knowledge of the genomic variation of the immunoglobulin loci is scarce. Here, we report 25 novel germline *IGHV* alleles as inferred from rearranged naïve B cell cDNA repertoires of 98 individuals. Thirteen novel alleles were selected for validation, out of which ten were successfully confirmed by targeted amplification and Sanger sequencing of non-B cell DNA. Moreover, we detected a high degree of variability upstream of the V-region in the 5'UTR, leader 1, and leader 2 sequences, and found that identical V-region alleles can differ in upstream sequences. Thus, we have identified a large genetic variation not only in the V-region but also in the upstream sequences of *IGHV* genes. Our findings challenge current approaches used for annotating immunoglobulin repertoire sequencing data.

Immunoglobulins are an important part of the adaptive immune system. They exert their function either as the antigen receptor of B cells that is essential for the antigen presentation capacity of these cells, or as secreted antibodies that survey extracellular fluids of the body. Immunoglobulins can bind a plethora of antigen epitopes via their paratopes, which are composed of combinations of heavy and light chain's variable regions. A huge diversity of paratopes is established by recombination of variable (V), diversity (D) (not in light chains) and joining (J) genes, and the pairing of heavy and light chains[1]. There is a large number of V, D, and J genes present on the heavy chain locus (chromosome 14, 14q32.33)[2] as well as the two light chain loci kappa (chromosome 2, 2p11.2) and lambda (chromosome 22, 22q11.2)[3].

35  These loci remain incompletely characterized due to the fact that they contain many repetitive

36  sequence segments with many duplicated genes[4], which makes it difficult to correctly assemble short

37  reads from whole genome sequencing. Single nucleotide polymorphisms as well as copy number

38  variations are in linkage disequilibrium and make up distinct haplotypes[4]. To this date, a limited

39  number of genomically sequenced [5-7] and inferred [8,9] haplotypes of the heavy chain and the two light

40  chain loci have been described. Different databases exist for genomic immune receptor DNA

41  sequences (IMGT/GENE-DB[10]), putative novel variants from inferred data (IgPdb[11]) or entire immune

42  receptor repertoires (OGRDB[12]).

43  The usage of immunoglobulin heavy chain variable (*IGHV)* genes and their mutational status are most

44  frequently studied in relation to cancer[13,14], responses to vaccines[15,16], or in autoimmune diseases[17-19].

45  Most *IGHV* genes have several allelic variants and more alleles are being discovered as a result of

46  adaptive immune receptor repertoire-sequencing (AIRR-seq)[20,21]. Software tools such as TIgGER[22,23],

47  IgDiscover[24] and partis[25] allow to infer germline alleles from such repertoire data. Based on these

48  inferred alleles, the data can then be input to other tools that infer haplotypes and repertoire

49  deletions[26]. Incorrect annotation could possibly lead to inferring wrong deletions and biased

50  assessments. Therefore, having a full overview of germline variants is essential for studying the
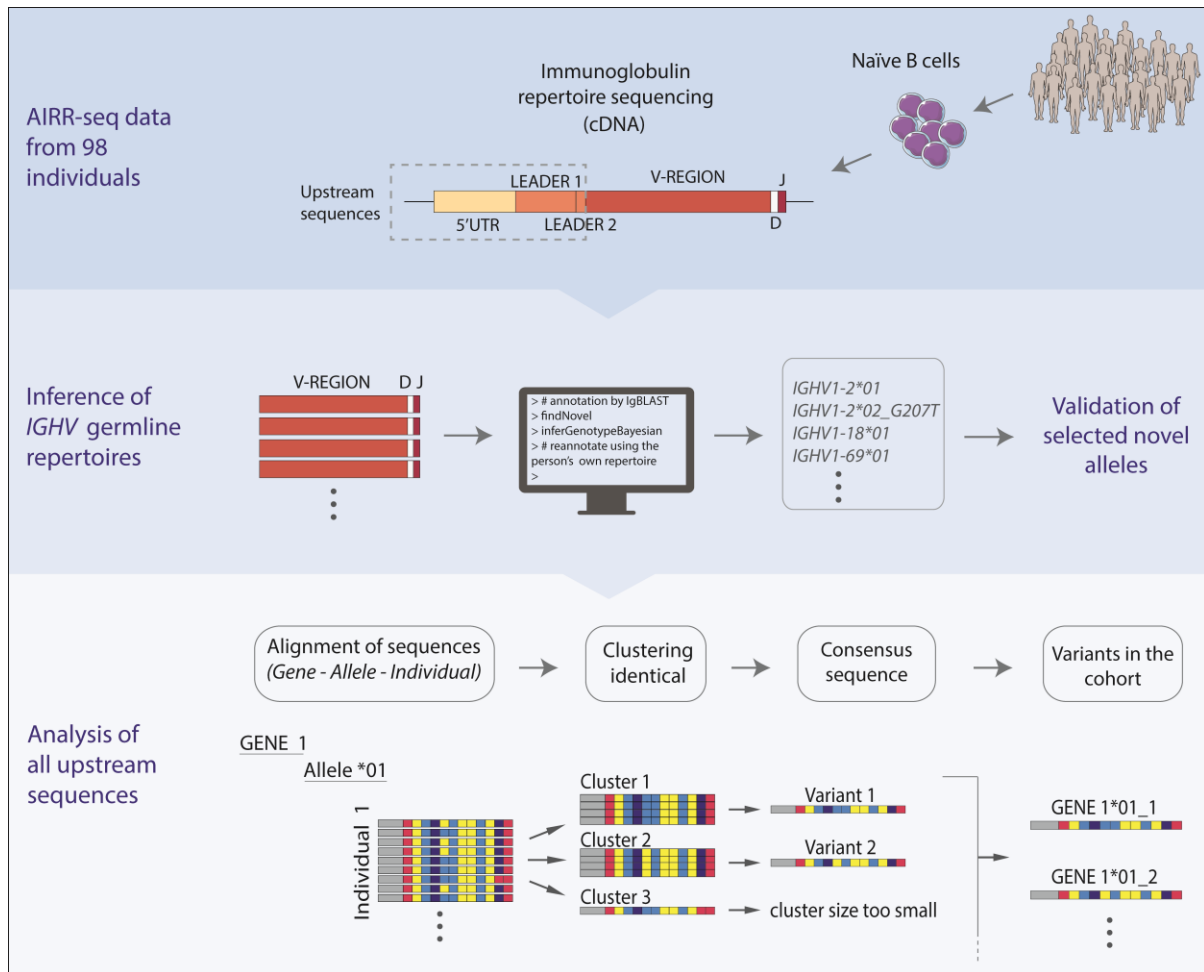
51  adaptive immune response with high accuracy.

52  Some allelic variants have been associated with increased disease susceptibility[27,28], yet the impact of

53  immunoglobulin gene variation on disease risks is still unknown[29]. These regions have not been

54  sufficiently covered in the numerous genome wide association studies performed to date. More

55  comprehensive maps of polymorphisms are required for proper analysis.

56  Here, we have used previously generated AIRR-seq data[30] from naïve B-cells of 98 Norwegian

57  individuals to identify novel *IGHV* alleles, a selection of which we then validated from genomic DNA

58  (gDNA) of non-B cells, i.e. T cells and monocytes. We also analyzed the sequences upstream of the

59  V-region, and constructed consensus sequences for the upstream variants present in the cohort.

60  These results expand our knowledge of this important locus and deepen our understanding of allelic

61  diversity within the Caucasian population. In addition, the result of this study can be used to improve

62  the accuracy of currently used bioinformatics tools for the analysis of immunoglobulin repertoire

63  sequencing data.

64

65  **RESULTS**

66  In this study, we used an AIRR-seq dataset from a cohort of 98 individuals[30] aiming to characterize

67  novel *IGHV* alleles that might be present in the data, as well as analyze the sequences upstream of

68  the V-region and create a table of previously unexplored upstream variants (Fig.1). To validate our

69  inferences from the AIRR-seq data analysis, genomic DNA of the same individuals was isolated from

70  non-B cells, i.e. T cells and monocytes. The reason for using non-B cells for validation was to avoid

71  capturing sequences with somatic hypermutation that occurs in primed B cells.

72

**Figure 1. Schematic representation of the data analysis.** In this study, we used material from a Norwegian cohort of 98 individuals[30]. Following the initial preprocessing of the data, we inferred the germline V-gene repertoires of all individuals in the cohort and identified novel alleles using the software suites TIgGER and IgDiscover. The availability of genomic DNA of the same individuals allowed us to verify some of our findings from the analysis of the AIRR-seq data. Since the validation attempts revealed polymorphisms outside of the V-region, we decided to analyze the upstream sequences, i.e. 5'UTR, leader 1 and leader 2. We used a custom approach for this analysis based on clustering identical variants. More details about the protocols and analysis can be found in the methods section.

We used two germline inference tools, TIgGER[22,23] and IgDiscover[24], to characterize novel alleles and to create a personalized germline reference of *IGHV* alleles for each individual (aka genotype). The purpose of using two different software tools was to increase our confidence in the inference of novel alleles. This study does not aim to serve as a comparison of the available software tools.

To increase the overlap between the different software results and to allow the discovery of novel alleles in genes with low expression, we adjusted selected TIgGER parameters, while keeping the IgDiscover parameters as default. Suspected false positive signals were filtered out from the novel allele candidates using mismatch frequency as described in Methods. The mismatch frequencies are depicted in Supplementary Fig.1. Novel allele candidates that were determined to be false positives contained mutations A152G, T154G and A85C (Supplementary Fig.1).

91

**92**     **Analysis of the V-region reveals 25 novel *IGHV* alleles**
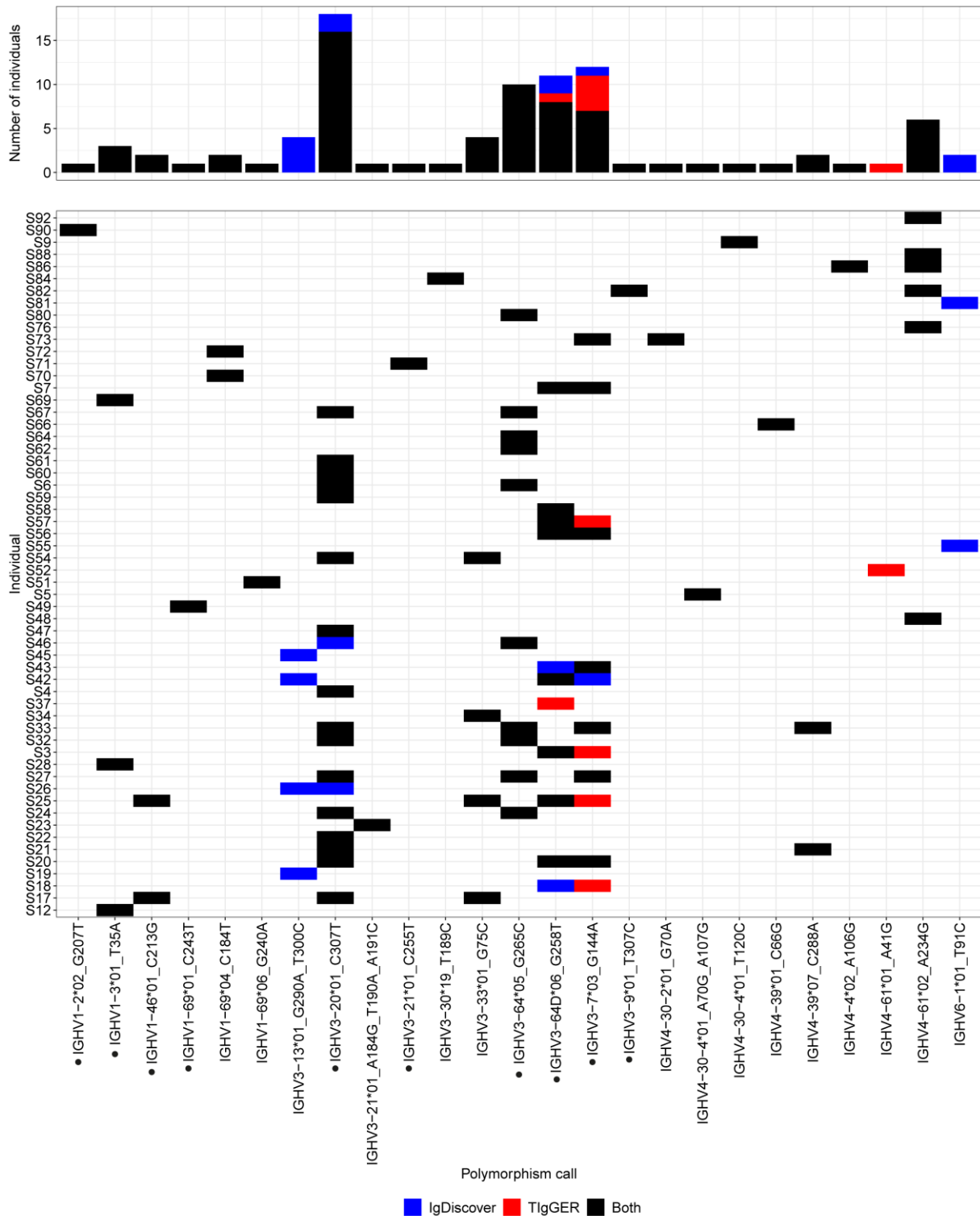
93     We first analyzed the usage of all genes and the different alleles carried by individuals in the cohort.

94     The relative usage of certain genes appeared to be strongly affected by the alleles present in the

95     inferred genotype. This was true for *IGHV1-2, IGHV1-46, IGHV3-11, IGHV3-43, IGHV3-48, IGHV3-53,*

96     *IGHV4-61,* and *IGHV5-51* (Supplementary Fig.2). Overview of the usage of all genes across all

97     individuals can be found in Supplementary Fig.2-3.

98     We inferred altogether 25 novel alleles (Fig.2), and we named them using the closest reference allele.

99     The majority of the novel alleles (22) were identified both with TIgGER and IgDiscover. In addition to

100     these, two novel alleles were found exclusively by IgDiscover, and one novel allele was found

101     exclusively by TIgGER.

102     Thirteen novel alleles were selected for validation by targeted amplification and subsequent Sanger

103     sequencing of gDNA (Supplementary Fig.4) of non-B cells, i.e.T cells and monocytes isolated by

104     fluorescence-activated cell sorting[30]. The validation primers are specified in the Supplementary Table

105     1. Out of those thirteen alleles, ten were successfully validated. These include *IGHV1-2*02_G207T,*

106     *IGHV1-3*01_T35A, IGHV1-46*01_C213G, IGHV1-69*01_C243T, IGHV3-7*03_G144A, IGHV3-*

107     *9*01_T307C, IGHV3-20*01_C307T, IGHV3-21*01_C255T, IGHV3-64*05_G265C,* and *IGHV3-*

108     *64D*06_G258T*. Surprisingly, *IGHV3-64*05_G265C* was found to originate from *IGHV3-64D* (Fig.6c).

109     Two of the novel alleles, namely *IGHV1-46*01_C213G* and *IGHV3-20*01_C307T,* have been recently

110     added to the IMGT database as *IGHV1-46*04* and *IGHV3-20*04* respectively.

111     Validation of the novel alleles revealed additional polymorphisms outside of the V-region. The allele

112     *IGHV3-64*06_G258T* has a polymorphism in leader 1 (position -21) in addition to the V-region

113     polymorphism. Genomic validation of *IGHV3-7*03_G144A* revealed a further polymorphism in the

114     intron. During validation of this allele, we also managed to amplify the genomic sequence of *IGHV3-*

115     *7*02*, which carried the previously reported polymorphism A318G[31]. This polymorphism was not

116     inferred from the AIRR-seq data in our study, since the default parameters of the inference tools are

117     set to detect polymorphisms up to position 312.

118     Attempts to validate *IGHV4-39*07_C288A*, *IGHV4-61*02_A234G*, and *IGHV6-1*01_T91C* were

119     unsuccessful. The gene-specific primers that were used for validation were designed based on the

120     current reference genome. However, the efficiency of the *IGHV4* primers was inferior, and Sanger

121     sequencing only revealed allele *01* of each gene, even in clearly heterozygous individuals.

122

**Figure 2. Novel *IGHV* alleles.** The software suites TIgGER and IgDiscover were used to infer a personal *IGHV* genotype for each individual and to infer previously undiscovered alleles. All novel alleles that are part of a genotype inferred by at least one of the methods appear on the x-axis. Alleles that were validated by Sanger sequencing are marked with a dot. Individuals with at least one novel allele lie on the y-axis and are labelled by their subject name. For each allele, the color of a tile (or a bar) represents the method of detection and genotype inference. The height of each bar on top represents the number of individuals for whom a certain allele was inferred and is part of a genotype.

130  **Analysis of upstream sequences yields a more complete and accurate germline reference**
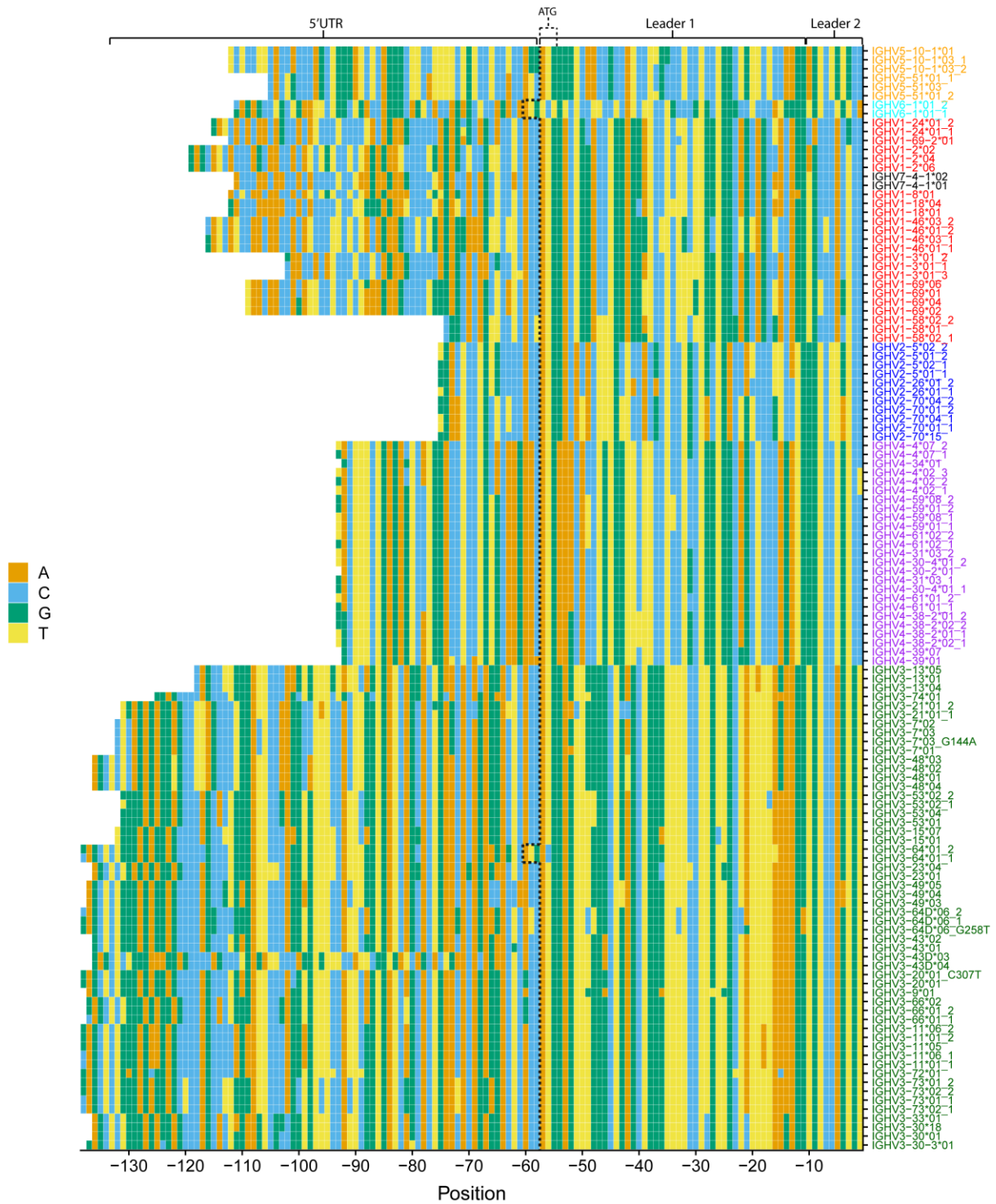131  **dataset**

132  As some of the validated novel alleles had additional polymorphisms in the intron or the leader
133  sequence, we extended our analysis of the AIRR-seq data beyond the V-region. Although introns are
134  not present in the AIRR-seq data, the sequences of the 5' untranslated region (5'UTR), leader 1, and
135  leader 2 lie upstream of the V-region and are present in the data thanks to the library preparation
136  method (Fig.1). We will refer to 5'UTR, leader 1, and leader 2 collectively as upstream sequences.

137  We decided to use the genotyped AIRR-seq data to characterize upstream sequence variants for all
138  genes and alleles. To extract the upstream sequences, we removed the VDJ and constant regions,
139  while keeping the original sequence's V-region annotation. Sequences from each individual were
140  processed separately. We observed slight variations in the length of 5'UTRs assigned to the same
141  gene. It is important to have matching length for clustering, as different lengths could mean that
142  identical sequences would not cluster together. To overcome this issue, for each gene we trimmed
143  the ends of 5' ends of the upstream sequences to match the most frequent length. We then took the
144  trimmed upstream sequences with the same allele annotation and clustered them. Each cluster of a
145  sufficient size gave rise to one consensus upstream sequence. This process was repeated for all
146  genes and alleles across all individuals. Finally, consensus sequences from all individuals were
147  combined to create an upstream germline reference dataset of the cohort (Fig.3). The number of
148  individuals carrying each of the variants is shown in Supplementary Fig.5.

149  According to the constructed germline reference dataset, the lengths of leader 1 (45 nt) and leader 2
150  (10 nt) sequences appear to be well conserved across most genes, with the exception of *IGHV3-*
151  *64*01* and *IGHV6-1*01* (Fig.3). The leader 1 sequences of these two genes are 3 nt longer, which
152  makes the position of ATG appear to be shifted upstream. The length of the 5'UTR is relatively
153  conserved within the same gene family, however, there is a large variability across different families.
154  Genes of the *IGHV2* family have the shortest 5'UTR, while the 5'UTRs of *IGHV3* genes are the
155  longest.

156  Comparison of  the consensus sequences in the cohort with the reference sequences obtained from
157  the IMGT/GENE-DB[10] revealed some discrepancies between our data and the reference database.
158  For example, the IMGT reference sequence of the allele *IGHV5-51*01* has T at position -3 in leader 2,
159  while the reference sequences of the other reference alleles have G at this position. However, in our
160  data, all *IGHV5-51* alleles have G at position -3, as illustrated in Fig.3. Our observation of G at
161  position -3 in *IGHV5-51*01* was validated by targeted amplification and Sanger-sequencing of *IGHV5-*
162  *51*01* from a homozygous individual (Supplementary Fig.6).

163

164

**Figure 3. Upstream germline reference dataset.** For each allele, consensus upstream sequences were built. Consensus sequences constructed from clusters with less than 10 sequences or with relative frequency < 0.1 were excluded. Each row represents a consensus upstream sequence of a V allele with 5' to 3' orientation. The colors of the tiles represent the different nucleotides. The coordinates on the x-axis describe the position of each nucleotide relative to the start of the V-region (5' to 3') and are therefore labeled as negative numbers. Alleles with more than one consensus sequence are marked with the allele name followed by an underscore and the respective consensus sequence number. For example, the two different consensus sequences for allele *IGHV3-64*01* are marked as *IGHV3-64*01_1* and *IGHV3-64*01_2*. The number of individuals who carry each variant are shown in Supplementary Fig.5.

174 **Length of 5'UTRs correlates with the distance between TATA-box and start codon**

175 As depicted in Fig.3, the length of the 5'UTR differs between *IGHV* gene families, but is relatively

176 conserved within a gene family. To investigate whether the different length of 5'UTRs among the

177 different families had any correlation with the distance from the promoter elements, we decided to

178 inspect the reference gDNA sequence from the IMGT database. We collected the available germline

179 reference sequences of the upstream flanking regions of V-gene promoters from the IMGT/GENE-DB

180 and aligned them to look for conserved patterns.

181 Using the sequences from the IMGT reference database, we determined the distance between the

182 ATG start codon and the reference or putative TATA-box. We found that this distance varied greatly

183 between different gene families. By comparing this distance to the 5'UTR length from the AIRR-seq

184 data, we observed that the distance between the ATG and the TATA-box correlated with the length of

185 the 5'UTR (Supplementary Fig.7). Sequences with longer ATG to TATA-box distance had longer
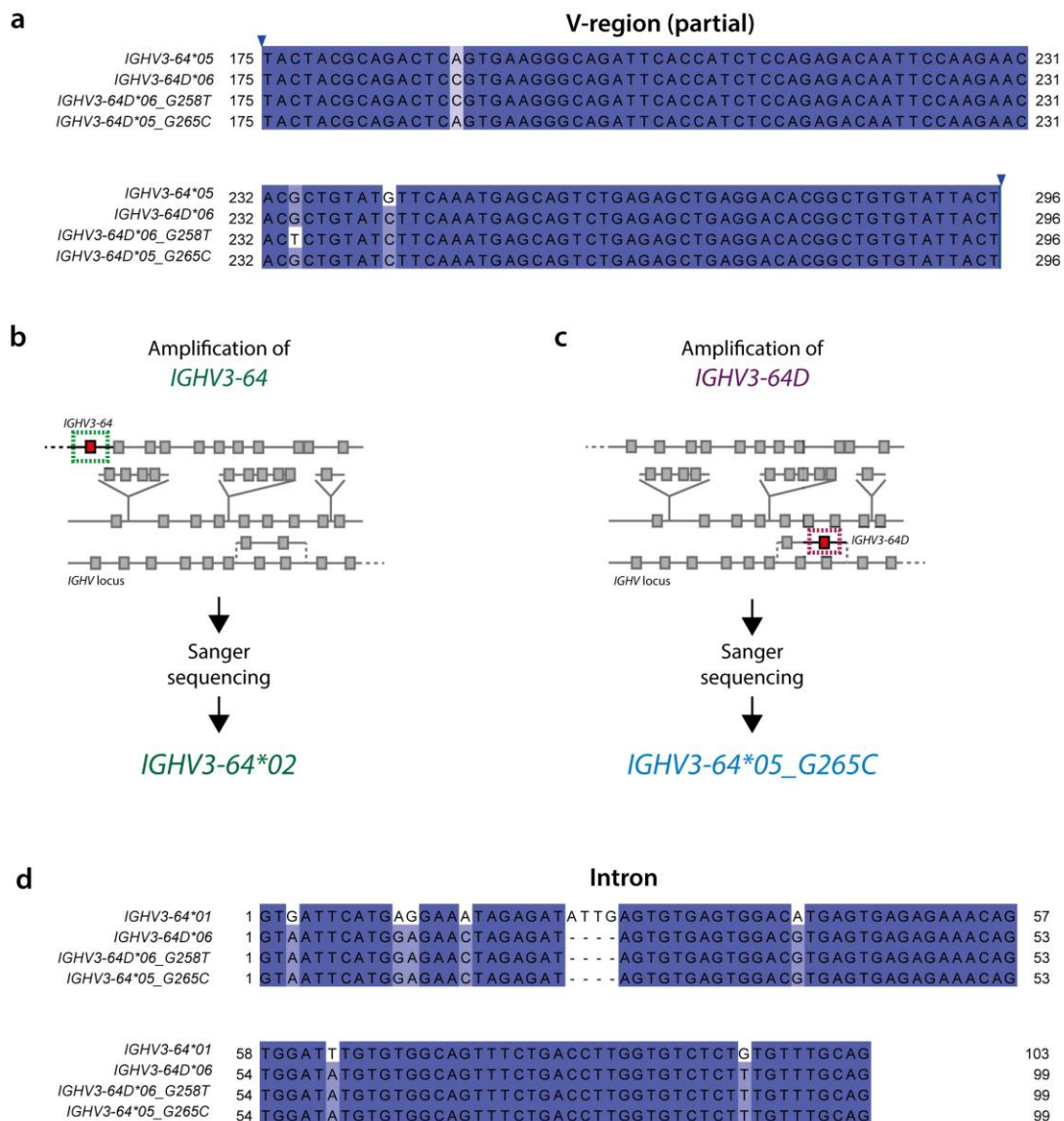
186 5'UTRs.

187 **Differences in the upstream sequences can aid allele annotation**

188 The novel allele *IGHV3-64*05_G265C* was initially not validated by amplification of the gene *IGHV3-*

189 *64,* as Sanger sequencing revealed only *IGHV3-64*02* in a selected individual carrying the suspected

190 polymorphism, and with no sequence corresponding to allele *05* being present (Fig.4b). Originally,

191 this allele was ambiguously annotated as deriving from either *IGHV3-64*05* or *IGHV3-64D*06*, as it

192 differs by one nucleotide from each of these alleles (Fig.4a).

193 The upstream sequences of *IGHV3-64* and *IGHV3-64D* differ all across their length, including the

194 5'UTR, leader 1, and leader 2 (Fig.3).  The upstream regions of the novel allele *IGHV3-64*05_G265C*

195 are identical to those of *IGHV3-64D*, which indicated that this is indeed an allele of *IGHV3-64D* and

196 not *IGHV3-64.* Therefore, we decided to amplify the gene *IGHV3-64D* using primers specific to the

197 duplicated gene only. This resulted in the novel allele being finally validated (Fig.4c). Upon obtaining

198 the full germline sequence of the novel allele, we observed that its intron matched the one of *IGHV3-*

199 *64D* and not *IGHV3-64* (Fig.4d).

200 The genes *IGHV3-43* and *IGHV3-43D* are another example of duplicated genes with differences in

201 the upstream sequences. Unlike the previous example, *IGHV3-43* and *IGHV3-43D* seem to have

202 identical leader 1 and leader 2 sequences but differ in the 5'UTR (Fig.3). However, not only genes,

203 but also some alleles of the same gene can be distinguished by their upstream sequences. The novel

204 allele *IGHV3-64D*06_G258T* differs from *IGHV3-64D*06* in one position located in leader 1. Similarly,

205 *IGHV4-39*01* and *IGHV4-39*07* have three differences within the 5'UTR; and the alleles *IGHV3-*

206 *43*01* and *02* differ in one position within the 5'UTR.

207

8

208

**Figure 4. Genomic validation of *IGHV3-64(D)* alleles.** (a) Alleles *IGHV3-64*05* and *IGHV3-64D*06* differ in only two positions within the V-region. To validate novel *IGHV3-64* and *IGHV3-64D* alleles found in the AIRR-seq data and ensure their correct annotation, we PCR amplified the genes *IGHV3-64* and *IGHV3-64D* from gDNA of selected individuals using gene-specific primers. (b, c) The process of validation of *IGHV3-64*05 G265C* depicted with a schematic *IGHV* locus representation. The novel allele was originally assigned as being closest to *IGHV3-64*05*, however, this allele was not amplified by primers specific for *IGHV3-64*. The novel allele was detected when *IGHV3-64D* was amplified. (d) Comparison of the intronic regions of *IGHV3-64*01*, *IGHV3-64D*06* and the novel *IGHV3-64D* alleles. The reference sequence of *IGHV3-64*05* in the IMGT database is partial and lacking the intron, and therefore could not be compared. The intron of the novel allele originally annotated as *IGHV3-64*05 G265C* matches the one of *IGHV3-64D*. The numbers in the alignments (a,d) do not follow the unique IMGT numbering.

220

221

222 **DISCUSSION**

223 Our analysis of the naïve B cell immunoglobulin repertoire data from 98 individuals revealed several

224 novel polymorphisms both in the coding and in the upstream sequences of *IGHV* genes. To our

225 knowledge, we are the first to provide a comprehensive overview of upstream (5'UTR, leader 1, and

226 leader 2) *IGHV* sequence variants in an AIRR-seq dataset. We managed to validate a number of

227 novel alleles by targeted amplification of genomic DNA of the same individuals. In addition, we report

228 the presence of G at position 318 instead of A in the gDNA sequence of *IGHV3-7*02*, which supports

229 the findings of previous studies[31,32].

230 We faced several issues with missing or incomplete genomic reference sequences, which

231 complicated the design of efficient primers for verification of novel alleles. Some of our validation

232 attempts were unsuccessful resulting only in the amplification of a "wild-type" allele without a

233 polymorphism. We suspect this might be caused by allelic dropout[33,34]. As we show in our upstream

234 sequence overview (Fig.4), alleles *IGHV4-39*01* and *IGHV4-39*07* differ at multiple positions within

235 the 5'UTRs. Our primers were designed to bind flanking sequences of the gene, and their design was

236 based on the current reference genome, which contains the allele *01* of *IGHV4-39.* Potential

237 differences in the primer binding regions could be the cause of a failure to amplify the novel alleles, in

238 this case *IGHV4-39*07_C288A*.

239 Although AIRR-seq studies are very useful for characterizing variation in immunoglobulin genes, one

240 of the main limitations are issues with gene and allele annotation[35]. The V-region is annotated based

241 on the most similar allele in the reference database. However, since the V genes are highly similar,

242 this annotation might not always be correct. Incorrect gene assignment could lead to potential

243 downstream errors in analysis. In our study, the novel allele originally annotated as *IGHV3-*

244 *64*05_G265C* was later found to be derived from the gene *IGHV3-64D*, located on a different part of

245 the *IGHV* locus than *IGHV3-64*. As previously shown[4,5,9] , *IGHV3-64D* is likely a part of an alternative

246 haplotype, since it was found to be deleted in many individuals, even in this cohort[30]. These two

247 genes differ in their upstream sequences, and thanks to this distinction, we were able to correctly

248 assign the novel allele to *IGHV3-64D* and validate it from gDNA.

249 Our results demonstrate that polymorphisms in the upstream regions can be utilized to improve

250 annotation methods presently employed. Having said that, the genetic variation in the sequences

251 upstream of the V-region is currently poorly characterized. Many reference sequences, which were

252 deposited to the IMGT germline database are partial and contain only the V-region sequence. It is

253 surprising that the genetic variation in the upstream regions is overlooked, considering the fact that

254 the leader regions are frequently used as primer binding sites for immunoglobulin repertoire library

255 preparation protocols[32,36,37].

256 The reason for the existence of upstream polymorphisms is unclear, but conceivably such

257 polymorphisms might have functional relevance by influencing stability of the mRNA or by affecting

258 the binding of regulatory proteins[38,39]. Further studies are needed to explore polymorphisms in the

10

259  upstream sequences and to determine whether they have any functional effect. Association of these

260  allelic variants with disease can be studied in sufficiently powered studies. In addition, more genomic

261  studies could be performed to characterize their promoters and other regulatory elements, which

262  might help explain the differences in expression levels across individuals.

263

264  **METHODS**

265  **AIRR Sequencing of naïve B-cells**

266  The data was obtained as a part of a previously published study[30]. In summary, naïve B cells from

267  100 individuals were sorted from peripheral blood mononuclear cells (PBMCs). The RNA was isolated

268  and quality checked before being sent to AbVitro, Inc for library preparation and sequencing on

269  Illumina MiSeq (2x300bp). About half of the cohort are celiac disease patients, and these subjects

270  were included to increase the diversity of the cohort. Of note, this study was not designed and

271  powered to perform comparative analysis of allelic frequencies between patients and controls.

272  **Amplification of target genomic regions**

273  Genomic DNA (gDNA) was isolated from previously sorted T cells and monocytes (CD19-

274  CD3+/CD14+)[30] using the QiaAmp DNA mini kit (Qiagen), and the concentration was measured on

275  Nanodrop.

276  Primers for validation were designed by PrimerBLAST using the reference genome as a template.

277  The nucleotide sequences of primers with additional details can be found in the Supplementary

278  material. For amplification of genes *IGHV3-7, IGHV3-20,* and *IGHV3-21,* primers from a recently-

279  published study [32] were used. All oligos were synthesized and purified (RP-cartridge) by Eurogentec.

280  The target regions of the gDNA were amplified by touch-down PCR using Q5® Hot Start High-Fidelity

281  DNA Polymerase (NEB). Approx. 100 - 200 ng gDNA from an individual with a suspected

282  polymorphism was used as a template. The PCR started with two cycles with the annealing

283  temperature of 70°C. The touch-down part of the PCR consisted of 10 cycles with the annealing

284  temperature decreasing from 70°C to 60°C by 1°C every cycle.  In the next 13 cycles, the annealing

285  temperature remained constantly at 60°C, and the last step of the PCR was the final extension at

286  72°C. The length of the PCR product varied depending on the amplified gene, ranging between 750bp

287  and 986bp.

288  **Cloning**

289  The PCR products were cleaned using the Monarch® DNA Gel Extraction Kit (NEB), and 3' end A-

290  overhangs were added by NEBNext® dA-Tailing Module (NEB). The A-tailed products were

291  subsequently cloned into pGEM®-T Easy vector (Promega) using the manufacturer's protocol. For

292  transformation, 4 μl of the ligation reaction were used to transform 90 μl XL10 $CaCl_2$-competent cells.

293  After transformation, 100 μl cells were plated on $LB_{amp}$ 50 μg/ml plates that have been previously

294  coated with IPTG/X-Gal (40 μl 100 mM IPTG + 16μl 50 mg/ml X-Gal). The IPTG/X-Gal treatment

295 allows for selection of successfully transformed colonies based on color. After overnight incubation at

296 37°C, white colonies were picked and the plasmids were isolated using the Monarch® Plasmid

297 Miniprep Kit (NEB). To verify that the picked colonies contain an insert of the correct size, a PCR was

298 performed using the same primers as for the amplification of gDNA, and the products were analyzed

299 by gel electrophoresis (1% agarose, 100 V, 35 min). The size of the PCR product was between 750-

300 986bp, depending on the gene amplified.

**Sanger sequencing**

302 Sanger sequencing of the plasmid DNA containing the correct-sized insert was performed by Eurofins.

303 The resulting sequences were trimmed to remove the vector and primer sequences. V-gene

304 annotation was done by IMGT/HighV-QUEST [40]. To check for polymorphisms in the introns, leader

305 regions and 5'UTRs, the trimmed sequences were aligned by MUSCLE [41,42] to the reference alleles of

306 the amplified gene, where available, and checked for polymorphisms. Alignments were visually

307 inspected in Jalview[43] and/or UGENE[44].

308 The sequences were named based on the amplified gene, followed by the closest reference allele

309 and the V-region polymorphism, which was determined by IMGT V-Quest [45] or by manual annotation

310 (in cases of ambiguous annotation).

311 The gDNA sequences of validated novel alleles were submitted to GenBank and subsequently to

312 IMGT.

**AIRR-seq data pre-processing**

314 The AIRR-seq data was pre-processed as described originally[30] using pRESTO [46]. Two individuals

315 were excluded from the analysis due to low sequencing depth (<2000).

**Novel allele discovery and genotype inference**

317 Genotype inference and novel allele discovery was also performed by TIgGER v 0.3.1 and IgDiscover

318 v0.11. The pre-processed sequences were annotated by IgBLAST 1.14.0[47] with modified parameters,

319 and the IMGT germline database (24) from January 2019 was used as a reference. The results of

320 alignment and genotype inference by TIgGER were processed using a similar pipeline to the one

321 used in http://www.vdjbase.org with slight modifications.

322 We experienced that the default settings resulted in incorrect annotation for some genes. This was

323 particularly obvious for the allele *IGHV5-51*03*, which was incorrectly annotated as *IGHV5-51*01* with

324 one mutation C45G, corresponding to the already known allele *03*. These two alleles differ only by

325 one nucleotide, and it was the length of the reference allele that seemed to affect whether or not the

326 sequence was correctly annotated by IgBLAST. The reference for *03* is 2 nt shorter than the

327 reference sequence for *01*, while sequences in our data corresponding to *IGHV5-51*03* were

328 matching the length of allele *01*. Adjusting the IgBLAST parameters --reward to 0 and --penalty to -3

329 resolved this annotation problem. These parameters were also induced manually in IgDiscover

330 alignment step.

331    For novel allele detection we tested the parameters of the TIgGER function "findNovelAlleles": 1)

332    germline_min to 50,100 and 200 (default). 2) j_max to 0.15 (default), 0.3 and 0.5. 3) min_seqs to 25

333    and 50 (default). Different parameters resulted in different sets of novel alleles identified. To allow for

334    discovery of novel alleles in lowly-expressed genes, we set the germline_min parameter to 50. The

335    rest of the parameters, including j_max and min_seqs, was left as default. The novel alleles were

336    further submitted for genotype inference, using a Bayesian approach, for each individual. As for

337    IgDiscover, the default parameters for novel allele and genotype calls were applied. Analysis of the

338    IgDiscover and TIgGER output was performed in R Studio version 3.6.0.

339    **Filtering out false positive suspects**

340    Errors that occur during the PCR reaction and/or sequencing could result in a false novel allele call.

341    To filter out the suspected false positive signals, we first determined the mismatch frequency for all

342    novel allele candidates. Novel allele candidates with low mismatch frequency were considered as

343    false positives. These included all alleles with mutation patterns A152G, T154G, and A85C. Although

344    the mismatch frequencies of sequences with the A85C polymorphism seemed to follow a bimodal

345    behavior (Supplementary Fig.1), the higher frequency mode that should correspond to heterozygous

346    individuals is centered around 20%, and not 50% as would be expected. As a result, they were not

347    considered as true novel alleles. On top of that, this polymorphism was only observed in four

348    individuals that were sequenced in a pilot separately from the other samples, and A to C mutation is

349    the most common substitution error in Illumina MiSeq[48].

350    **Analysis of gene and allele usage**

351    Following the inference of genotype for each individual, we used IgBLAST 1.14.0[47] to re-align each

352    individual's sequences with their own personalized germline *IGHV* database as inferred by TIgGER.

353    To compare the relative gene usage in individuals with different allele combination, we selected

354    sequences with V-region length >200 and up to 3 mutations. Since the duplicated genes *IGHV3-*

355    *23*01* and *IGHV3-23D*01*; *IGHV1-69*01* and *IGHV1-69D*01*; *IGHV2-70*04* and *IGHV2-70D*04* have

356    identical V-regions, they often result in ambiguous allele assignment. Annotation for sequences with

357    ambiguous allele assignments for these genes were renamed *IGHV3-23*01D*, *IGHV1-69*01D* and

358    *IGHV2-70*04D*, respectively. Additionally, *IGHV3-30-5*01* and *IGHV3-30*18* are also identical; and

359    we renamed them as *IGHV3-30X*doub*; and *IGHV3-30X*trip* if the sequence annotation also

360    contained *IGHV3-30*01* as another possible assignment. All remaining sequences with multiple allele

361    annotations were filtered out. To plot the relative gene usage, we first calculated the relative usage

362    fraction of each allele of a gene separately. Afterwards, we summed up the relative usage fractions of

363    alleles of the same gene and plotted the relative usage of each gene across all individuals.

364

365    **Inference of upstream sequences (5'UTR, leader 1 and leader 2)**

366    We decided to look at the upstream regions that consist of (5'-3') 5'UTR, leader 1, and leader 2. For

367    the analysis of the upstream regions, only sequences with up to 3 mutations in the V-region (after

368    novel allele inference and genotyping) and single assignment V-call were selected. For each

369    individual, the V-region sequences were trimmed away and the remaining upstream sequences of the

370    same V-gene were aligned by the last nucleotide of leader 2 sequence and flipped 3'-5'.

371    Since the length of the 5'UTR sequences of the same gene in AIRR-seq data can vary due to whole

372    VDJ sequence length and sequencing length limitations, we needed to determine where to trim the

373    longer sequences. To do this, we first filled the ends of sequences with Ns to match the length of the

374    longest sequence for the respective gene. We then trimmed all sequences after the first position, at

375    which 95% sequences contained N.

376    After that, for each allele and for each individual, we removed all artificially added Ns. Next, we

377    estimated sequence lengths, and lengths with frequency above 0.1 were considered frequent.

378    Sequences shorter than the shortest frequent sequence length were filtered out and sequences

379    longer than the longest frequent sequence length were trimmed to match its length. By applying

380    ClusterSets.py (--ident 0.999, --length 0.5) and BuildConsensus.py (--freq 0.6) from pRESTO, we

381    constructed clusters that resulted in consensus sequences for each allele. For each cluster we

382    calculated its frequency based on the number of sequences assigned to it. Clusters with frequency

383    below 0.1 or with less than 10 sequences were removed.

384    For each allele, consensus sequences from all individuals, were trimmed to match the shortest

385    consensus sequence, and identical sequences were re-collapsed by allele and individual. For some of

386    the consensus sequences, one of the nucleotides was marked with ambiguous assignment (N) by

387    BuildConsensus.py function. In such cases, the original cluster was split into two clusters based on

388    the ambiguous assignment and consensus sequences were reconstructed manually. Finally, to create

389    the consensus upstream sequences, for each allele the trimmed sequences were submitted to

390    ClusterSets.py (--ident 1.0, --length 1.0) and BuildConsensus.py (--freq 0.6) functions and as a result,

391    for each gene and allele a set of consensus V upstream sequences were gathered. In the last step,

392    we compared and collapsed identical sequences from all individuals to create a database of upstream

393    sequences in the cohort.

394    **Analysis of the reference germline upstream sequences**

395    Reference germline sequences of the upstream sequences, including the 5'UTR, were obtained from

396    the IMGT GENE-DB and by searching through the IMGT "Gene tables" in order to get an alternative

397    longer sequence if available. The reference upstream sequences longer than 150 nt were aligned

398    using the MUSCLE tool at EMBL-EBI [42], and the alignment was visualized by Jalview [43] to look for

399    conserved regions. The obtained consensus sequences of conserved regions were compared to

400    IMGT resources for annotation. The TATA-boxes were determined based on either the reference

401    annotation by IMGT, searching through previous studies, or by looking for a TA-rich region

402    downstream of the octamer. Promoters studied by older studies include that of *IGHV6*[49] (with two

403    TATA-boxes) and *IGHV1*[50].

404    *IGHV2* analysis is based on the available upstream reference sequences of *IGHV2-5*01,*02* and

405    *IGHV2-70D*04,*14. IGHV3* schematic promoter representation was based on the upstream reference

406     genomic sequences of *IGHV3-43\*01, IGHV3-48\*02, IGHV3-49\*03, IGHV3-64\*02, IGHV3-64D\*06* and

407     the genomic sequences obtained by Sanger sequencing of *IGHV3-7\*02* and *IGHV3-64D\*06*. The

408     IGHV4 schematic representation of the promoter was based the reference genomic sequences of

409     *IGHV4-4\*07* and *\*08*; *IGHV4-28\*01,\*02,\*07*; *IGHV4-30-2\*06*; *IGHV4-30-4\*07*; *IGHV4-31\*02*; *IGHV4-*

410     *34\*01,\*02,\*11; IGHV4-38-2\*02; IGHV4-39\*01; IGHV4-59\*01,\*02,\*11;* and *IGHV4-61\*01,\*08,\*09.*

411     **Data availability**

412     The pipeline for novel allele discovery and genotype processing using the software tools TIgGER and

413     IgBLAST is available on the VDJbase website (https://www.vdjbase.org). Custom code for the

414     analysis of upstream sequences is available at https://bitbucket.org/yaarilab/cluster_5utr/src/master/ .

415     Sanger sequences of validated *IGHV* alleles have been deposited in the GenBank under accession

416     numbers: MN337615 (*IGHV1-2\*02_G207T*), MN337616 (*IGHV1-3\*01_T35A*), MN337617 (*IGHV1-*

417     *46\*01_C213G*), MN337618 (*IGHV1-69\*01_C243T*), MN337619 (*IGHV3-7_G144A_T300C*),

418     MN337620 (*IGHV3-7\*02_A318G*), MN337621 (*IGHV3-9\*01_T307C*), MN337622 (*IGHV3-*

419     *20\*01_C307T*), MN337623 (*IGHV3-21\*01_C255T*), MN337624 (*IGHV3-64D\*06_G258T*), and

420     MN337625 (*IGHV3-64D\*06_C210A*).

421

422     **ACKNOWLEDGEMENTS**

437

438

439

**AUTHOR CONTRIBUTIONS**

L.M.S. and G.Y. conceived and supervised the project; I.M., I.L. and O.S. carried out the experimental work; M.G., I.M., A.P., and G.Y. analyzed the data; I.M., M.G., and L.M.S. wrote the paper. All authors edited the manuscript.

**COMPETING INTERESTS**

No conflict of interests declared.

**REFERENCES**

1.  Murphy, K. & Weaver, C. *Janeway's immunobiology*, (Garland Science, 2016).
2.  McBride, O.W. *et al.* Localization of human variable and constant region immunoglobulin heavy chain genes on subtelomeric band q32 of chromosome 14. *Nucleic Acids Research* **10**, 8155-8170 (1982).
3.  McBride , O. *et al.* Chromosomal location of human kappa and lambda immunoglobulin light chain constant region genes. *The Journal of Experimental Medicine* **155**, 1480-1490 (1982).
4.  Watson, C.T. & Breden, F. The immunoglobulin heavy chain locus: genetic variation, missing data, and implications for human disease. *Genes & Immunity* **13**, 363-373 (2012).
5.  Watson, Corey T. *et al.* Complete haplotype sequence of the human immunoglobulin heavy-chain variable, diversity, and joining genes and characterization of allelic and copy-number variation. *The American Journal of Human Genetics* **92**, 530-546 (2013).
6.  Matsuda, F. *et al.* The complete nucleotide sequence of the human immunoglobulin heavy chain variable region locus. *The Journal of Experimental Medicine* **188**, 2151 (1998).
7.  Watson, C.T. *et al.* Sequencing of the human IG light chain loci from a hydatidiform mole BAC library reveals locus-specific signatures of genetic diversity. *Genes & Immunity* **16**, 24-34 (2015).
8.  Kidd, M.J. *et al.* The inference of phased haplotypes for the immunoglobulin H chain V region gene loci by analysis of VDJ gene rearrangements. *The Journal of Immunology* **188**, 1333 (2012).
9.  Kirik, U., Greiff, L., Levander, F. & Ohlin, M. Parallel antibody germline gene and haplotype analyses support the validity of immunoglobulin germline gene inference and discovery. *Molecular Immunology* **87**, 12-22 (2017).
10. Giudicelli, V., Chaume, D. & Lefranc, M.-P. IMGT/GENE-DB: a comprehensive database for human and mouse immunoglobulin and T cell receptor genes. *Nucleic Acids Research* **33**, D256-D261 (2005).
11. The Immunoglobulin Polymorphism IgGRdb (IgPdb).
12. Lees, W. *et al.* OGRDB: a reference database of inferred immune receptor genes. *Nucleic Acids Research* **48**, D964-D970 (2019).
13. Rodríguez-Vicente, A.E. *et al.* Next-generation sequencing in chronic lymphocytic leukemia: recent findings and new horizons. *Oncotarget* **8**, 71234-71248 (2017).
14. Ghiotto, F. *et al.* Mutation pattern of paired immunoglobulin heavy and light variable domains in chronic lymphocytic leukemia B cells. *Molecular medicine (Cambridge, Mass.)* **17**, 1188-1195 (2011).
15. Wang, C. *et al.* B-cell repertoire responses to varicella-zoster vaccination in human identical twins. *Proceedings of the National Academy of Sciences* **112**, 500 (2015).

16

484   16.   Galson, J.D. *et al.* B-cell repertoire dynamics after sequential hepatitis B vaccination and
485         evidence for cross-reactive B-cell activation. *Genome Medicine* **8**, 68 (2016).
486   17.   Roy, B. *et al.* High-throughput single-cell analysis of B cell receptor usage among
487         autoantigen-specific plasma cells in celiac disease. *The Journal of Immunology* **199**, 782
488         (2017).
489   18.   Di Niro, R. *et al.* High abundance of plasma cells secreting transglutaminase 2–specific IgA
490         autoantibodies with limited somatic hypermutation in celiac disease intestinal lesions.
491         *Nature Medicine* **18**, 441 (2012).
492   19.   Bashford-Rogers, R.J.M., Smith, K.G.C. & Thomas, D.C. Antibody repertoire analysis in
493         polygenic autoimmune diseases. *Immunology* **155**, 3-17 (2018).
494   20.   Brown, A.J. *et al.* Augmenting adaptive immunity: progress and challenges in the
495         quantitative engineering and analysis of adaptive immune receptor repertoires. *Molecular*
496         *Systems Design & Engineering* **4**, 701-736 (2019).
497   21.   Yaari, G. & Kleinstein, S.H. Practical guidelines for B-cell receptor repertoire sequencing
498         analysis. *Genome Medicine* **7**, 121 (2015).
499   22.   Gadala-Maria, D. *et al.* Identification of subject-specific immunoglobulin alleles from
500         expressed repertoire sequencing data. *Frontiers in Immunology* **10**(2019).
501   23.   Gadala-Maria, D., Yaari, G., Uduman, M. & Kleinstein, S.H. Automated analysis of high-
502         throughput B-cell sequencing data reveals a high frequency of novel immunoglobulin V gene
503         segment alleles. *Proceedings of the National Academy of Sciences* **112**, E862 (2015).
504   24.   Corcoran, M.M. *et al.* Production of individualized V gene databases reveals high levels of
505         immunoglobulin genetic diversity. *Nature Communications* **7**, 13642 (2016).
506   25.   Ralph, D.K. & Matsen, F.A.I.V. Per-sample immunoglobulin germline inference from B cell
507         receptor deep sequencing data. *PLOS Computational Biology* **15**, e1007133 (2019).
508   26.   Peres, A., Gidoni, M., Polak, P. & Yaari, G. RAbHIT: R antibody haplotype inference tool.
509         *Bioinformatics* **35**, 4840-4842 (2019).
510   27.   Parks, T. *et al.* Association between a common immunoglobulin heavy chain allele and
511         rheumatic heart disease risk in Oceania. *Nature Communications* **8**, 14946 (2017).
512   28.   Avnir, Y. *et al.* IGHV1-69 polymorphism modulates anti-influenza antibody repertoires,
513         correlates with IGHV utilization shifts and varies by ethnicity. *Scientific Reports* **6**, 20842
514         (2016).
515   29.   Watson, C.T., Glanville, J. & Marasco, W.A. The Individual and Population Genetics of
516         Antibody Immunity. *Trends in Immunology* **38**, 459-470 (2017).
517   30.   Gidoni, M. *et al.* Mosaic deletion patterns of the human antibody heavy chain gene locus
518         shown by Bayesian haplotyping. *Nature Communications* **10**, 628 (2019).
519   31.   Thörnqvist, L. & Ohlin, M. Critical steps for computational inference of the 3′-end of novel
520         alleles of immunoglobulin heavy chain variable genes - illustrated by an allele of IGHV3-7.
521         *Molecular Immunology* **103**, 1-6 (2018).
522   32.   Vázquez Bernat, N. *et al.* High-quality library preparation for NGS-based immunoglobulin
523         germline gene inference and repertoire expression analysis. *Frontiers in Immunology*
524         **10**(2019).
525   33.   Blais, J. *et al.* Risk of misdiagnosis due to allele dropout and false-positive PCR artifacts in
526         molecular diagnostics: analysis of 30,769 genotypes. *The Journal of Molecular Diagnostics* **17**,
527         505-514 (2015).
528   34.   Soulsbury, C.D., Iossa, G., Edwards, K.J., Baker, P.J. & Harris, S. Allelic dropout from a high-
529         quality DNA source. *Conservation Genetics* **8**, 733-738 (2007).
530   35.   Smakaj, E. *et al.* Benchmarking immunoinformatic tools for the analysis of antibody
531         repertoire sequences. *Bioinformatics* (2019).
532   36.   René, C. *et al.* Comprehensive characterization of immunoglobulin gene rearrangements in
533         patients with chronic lymphocytic leukaemia. *Journal of cellular and molecular medicine* **18**,
534         979-990 (2014).

37. Vergani, S. *et al.* Novel method for high-throughput full-length IGHV-D-J sequencing of the immune repertoire from bulk B-cells with single-cell resolution. *Frontiers in Immunology* **8**(2017).

38. Steri, M., Idda, M.L., Whalen, M.B. & Orrù, V. Genetic variants in mRNA untranslated regions. *WIREs RNA* **9**, e1474 (2018).

39. Burke, T.W. & Kadonaga, J.T. The downstream core promoter element, DPE, is conserved fromDrosophila to humans and is recognized by TAFII60 of Drosophila. *Genes & Development* **11**, 3020-3031 (1997).

40. Alamyar, E., Duroux, P., Lefranc, M.-P. & Giudicelli, V. IMGT® Tools for the nucleotide analysis of immunoglobulin (IG) and T cell receptor (TR) V-(D)-J repertoires, polymorphisms, and IG mutations: IMGT/V-QUEST and IMGT/HighV-QUEST for NGS. in *Immunogenetics: Methods and Applications in Clinical Practice* (eds. Christiansen, F.T. & Tait, B.D.) 569-604 (Humana Press, Totowa, NJ, 2012).

41. Edgar, R.C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research* **32**, 1792-1797 (2004).

42. Madeira, F. *et al.* The EMBL-EBI search and sequence analysis tools APIs in 2019. *Nucleic acids research* **47**, W636-W641 (2019).

43. Waterhouse, A.M., Procter, J.B., Martin, D.M.A., Clamp, M. & Barton, G.J. Jalview Version 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189-1191 (2009).

44. Okonechnikov, K., Golosova, O., Fursov, M. & the, U.t. Unipro UGENE: a unified bioinformatics toolkit. *Bioinformatics* **28**, 1166-1167 (2012).

45. Giudicelli, V., Chaume, D. & Lefranc, M.-P. IMGT/V-QUEST, an integrated software program for immunoglobulin and T cell receptor V–J and V–D–J rearrangement analysis. *Nucleic Acids Research* **32**, W435-W440 (2004).

46. Vander Heiden, J.A. *et al.* pRESTO: a toolkit for processing high-throughput sequencing raw reads of lymphocyte receptor repertoires. *Bioinformatics* **30**, 1930-1932 (2014).

47. Ye, J., Ma, N., Madden, T.L. & Ostell, J.M. IgBLAST: an immunoglobulin variable domain sequence analysis tool. *Nucleic Acids Research* **41**, W34-W40 (2013).

48. Schirmer, M. *et al.* Insight into biases and sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Research* **43**, e37-e37 (2015).

49. Sun, Z. & Kitchingman, G.R. Bidirectional transcription from the human immunoglobulin VH6 gene promoter. *Nucleic Acids Research* **22**, 861-868 (1994).

50. Eaton, S. & Calame, K. Multiple DNA sequence elements are necessary for the function of an immunoglobulin heavy chain promoter. *Proceedings of the National Academy of Sciences* **84**, 7634 (1987).

# Polymorphisms in immunoglobulin heavy chain variable genes and their upstream regions

Ivana Mikocziova[1†*], Moriah Gidoni[2†], Ida Lindeman[1], Ayelet Peres[2], Omri Snir[1], Gur Yaari[2‡], Ludvig M. Sollid[1‡]

[1] K.G.Jebsen Centre for Celiac Disease Research and Department of Immunology, University of Oslo and Oslo University Hospital, 0372 Oslo, Norway
[2] Faculty of Engineering, Bar Ilan University, Ramat Gan 5290002, Israel
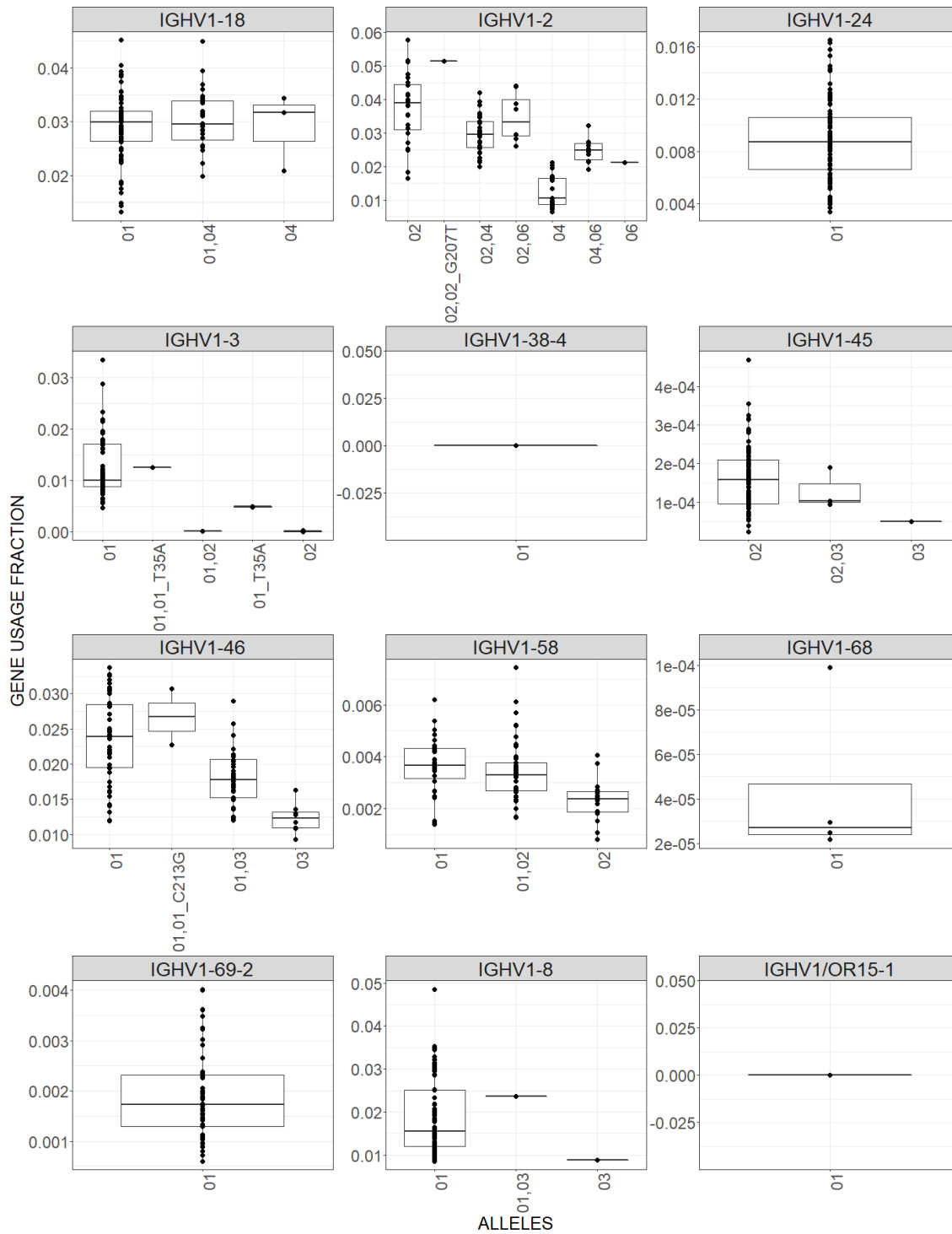
† Joint First Authors

‡ Joint Last Authors

* To whom correspondence should be addressed. Email: ivana.mikocziova@medisin.uio.no

# Supplementary Material

**Supplementary Figure 1. Germline mismatch frequency.** For each individual, relative frequencies of polymorphisms (y-axis) were calculated for positions in sequences aligned to an allele, for which a novel allele was inferred in the dataset (x-axis). Each dot represents a mismatch frequency for an individual for a certain allele and nucleotide. The color of the dot represents the nucleotide that does not match the germline.

**Supplementary Figure 2. Usage of genes across individuals in the cohort.** For each allele of a gene, we calculated its relative usage fraction in each individual. The usage fractions of alleles of the same gene in the same individual were then summed, revealing the gene's usage fraction. The x-axis shows the inferred allele, or multiple alleles, that were found in an individual's inferred genotype. Each dot represents one individual. The y-axis shows the relative usage fraction of a gene within the expressed repertoire. The bar represents the median value. A bias can be observed in some genes, where the median gene usage is higher in individuals homozygous for a specific allele than those homozygous for another allele. This figure continues on the next five pages.

Supplementary Figure 2. continued

**Supplementary Figure 2. continued**

**Supplementary Figure 2. continued**

**Supplementary Figure 2. continued**

**Supplementary Figure 2. Continued**

**IGHV1-69**

**Supplementary Figure 3. Usage of *IGHV1-69* across individuals in the cohort.** Relative usage fraction was calculated for each allele separately and in each individual, and the relative fractions of all expressed alleles were summed up. Different combinations of expressed alleles are shown on the x-axis, and the summed gene usage fraction is shown on the y-axis. Each dot represents one individual. The bar in the boxplot represents the median value.

**Supplementary Figure 4.** Sanger sequencing results. Ten novel alleles were validated by targeted amplification and subsequent Sanger sequencing. The trace files were aligned to reference sequences from IMGT GENE-DB[1] and visualised by UGENE[2]. This figure continues on the next page.
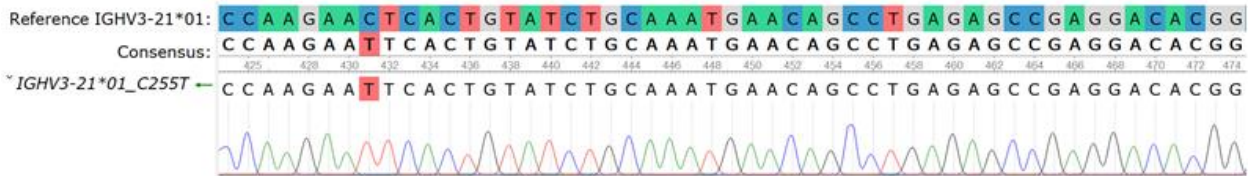
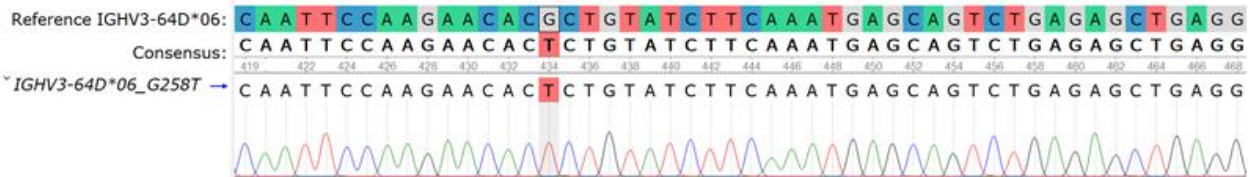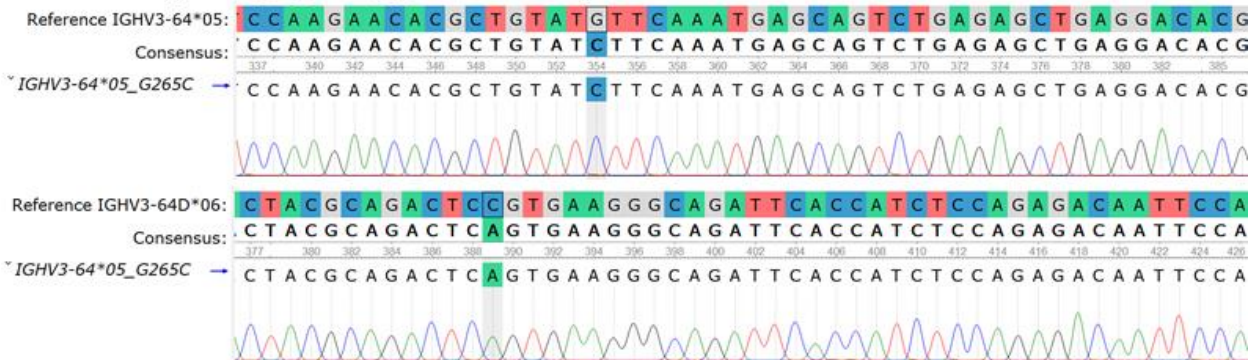**Supplementary Figure 4. continued**

**Supplementary Table 1. Primers used for genomic validation.**

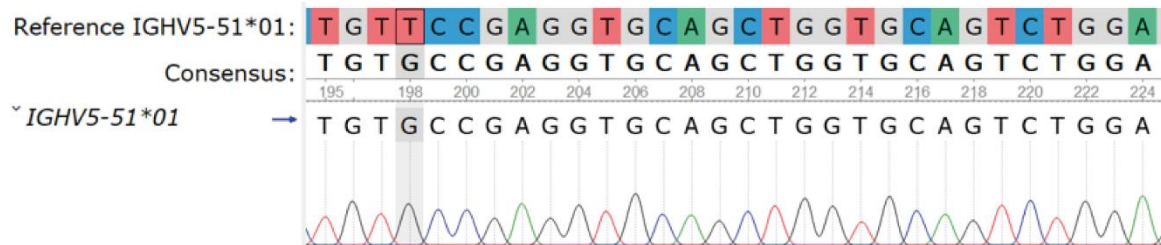| Primer name | Nt sequence (5'-3') |
|---|---|
| IGHV1-2_fwd | CGGGAACTTGTTTTCAGCAGAC |
| IGHV1-2_rev | TTTCATTTCTCAGCCCCAGCA |
| IGHV1-3_fwd | TCCAGTGGGAGAAGCTCTGT |
| IGHV1-3_rev | GTCATTTCCTCCATGCCAGC |
| IGHV1-46_fwd | CTGTGTGGCAGATGGGACAT |
| IGHV1-46_rev | TACTGAGTGTGGCCTTTCCC |
| IGHV1-69_fwd | TGGGAGCACAGCTCATCA |
| IGHV1-69_rev | CACTCTCAGGATGTGGGTTT |
| IGHV3-9_fwd | AGGACTCACCATGGAGTTGG |
| IGHV3-9_rev | TTTTTGTCTGGGCTCTCGCT |
| IGHV3-11_fwd | CAGCGTCCCACTAGAGCTTG |
| IGHV3-11_rev | CTGCAGGGAGGTTTGTGTCT |
| IGHV3-23_fwd | ATGCAAATAGAGCCCTCCGTCT |
| IGHV3-23_rev | TTCTGTCCCAGGACTGATTGCG |
| IGHV3-64_fwd | TTTGGGCTGAGCTGGGTTTT |
| IGHV3-64_rev | CAGGGAGGTTTCTGCATGGT |
| IGHV3-64D_fwd | AAGGACACTCTCATCTGCCC |
| IGHV3-64D_rev | CTCCTTGTGCACCTGCCTC |
| IGHV5-51_fwd | GAGAGGGACAATAGCAGGGTGTA |
| IGHV5-51_rev | CATATTGGAGAGGTGCCTGTTAGG |
| IGHV6-1_fwd | AGTCACCAGAGCTCCAGACA |
| IGHV6-1_rev | GCTCACACTGACTTCCCCTC |
| **Primers from Vázquez Bernat _et al_.[3]:** | |
| IGHV3-7R | CCTGGGGAAATTTGACGACGAGGCA |
| IGHV3-7F | GGGTACAGCCTATTCCTCCAGCA |
| IGHV3-20R | GCACCTGGTCCCTGAGTTTACTGTGTTC |
| IGHV3-20F | CACGGGCCAGACAGTGAGACTGG |
| IGHV3-21R | CGCCGCAGGCCATGACAGGAAGC |
| IGHV3-21F | CAGCGTCCCACCCTAGAGCTTGT |

**Supplementary Figure 5. Number of individuals carrying the consensus 5'UTR sequence.** Consensus 5'UTR sequences for each *IGHV* allele across all individuals were gathered and clustered to create a 5'UTR database. The length of each bar (x-axis) is the cluster size for a specific 5'UTR allele (y-axis), i.e. the number of individuals carrying the variant.
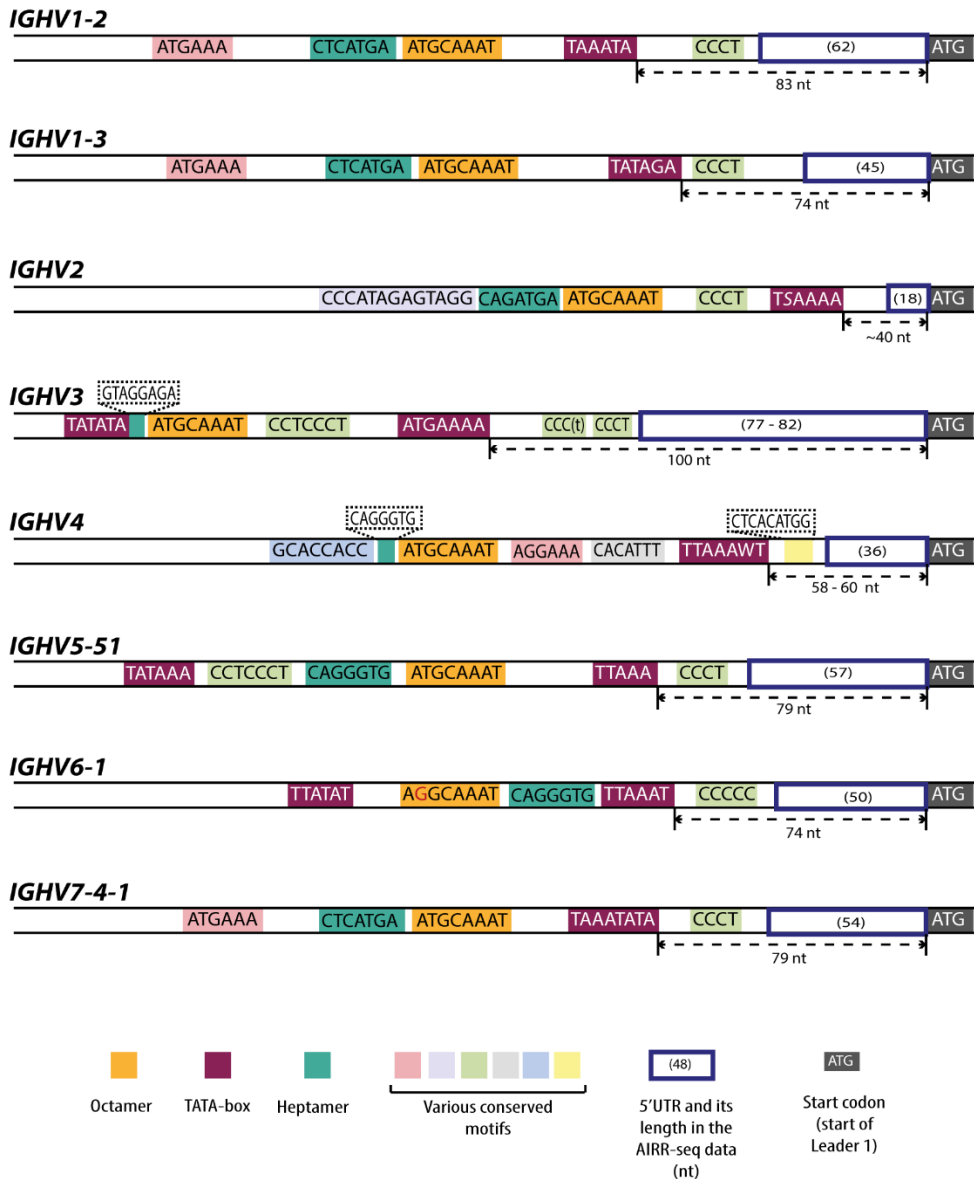
a



b



**Supplementary Figure 6. Validation of the _IGHV5-51*01_ leader 2.** The leader 2 part of _IGHV5-51*01_ in all individuals in our cohort differed from the reference in the IMGT database. (a) Alignment of the reference leader 2 sequences of selected _IGHV5-51_ alleles. (b) An individual from our cohort homozygous for _IGHV5-51*01_ was selected and _IGHV5-51_ was amplified using gene-specific primers (shown in Supplementary Table 1). Sanger sequencing of the amplified product revealed that the leader 2 of _IGHV5-51*01_ indeed differs from the reference.

**Supplementary Figure 7. Schematic representation of the *IGHV* promoter regions.** Reference upstream genomic sequences, including the promoter region were retrieved from the IMGT germline database and schematically depicted. Conserved motifs were identified by aligning all available 5'UTR and promoter reference sequences (> 150 nt) by MUSCLE and by searching for regions with high levels of homology. TATA-box sequences (in maroon) of some genes have been previously reported. For the remaining genes, we identified a putative TATA-box by searching for a TA-rich sequence. The octamer (in yellow) is well characterized and highly conserved across all genes. The heptamer (in dark turquoise) was only characterized for *IGHV1* genes. In the other genes, we identified putative heptamers by searching for a conserved sequence upstream of the octamer. Various conserved motifs with unknown function were also identified (pastel colors). The ATG start codon is shown in grey. The 5'UTRs that are found in the AIRR-seq data are lined in dark blue, and their typical length in the repertoire sequencing data is shown in brackets. The length of the 5'UTRs correlated with the distance between the ATG and the TATA-box.

15

**Supplementary References**

1.  Giudicelli, V., Chaume, D. & Lefranc, M.-P. IMGT/GENE-DB: a comprehensive database for human and mouse immunoglobulin and T cell receptor genes. *Nucleic Acids Research* **33**, D256-D261 (2005).
2.  Okonechnikov, K., Golosova, O., Fursov, M. & the, U.t. Unipro UGENE: a unified bioinformatics toolkit. *Bioinformatics* **28**, 1166-1167 (2012).
3.  Vázquez Bernat, N. *et al.* High-quality library preparation for NGS-based immunoglobulin germline gene inference and repertoire expression analysis. *Frontiers in Immunology* **10**(2019).