# A comparison of immunoglobulin IGHV, IGHD and IGHJ genes in wild-derived and classical inbred mouse strains. — Source link ⧉

Corey T. Watson, Justin T. Kos, William S. Gibson, Leah C. Newman ...+5 more authors

**Institutions:** University of Louisville, Icahn School of Medicine at Mount Sinai, German Cancer Research Center, Garvan Institute of Medical Research ...+1 more institutions

Related papers:

- A comparison of immunoglobulin IGHV, IGHD and IGHJ genes in wild-derived and classical inbred mouse strains

- The mouse antibody heavy chain repertoire is germline-focused and highly variable between inbred strains.

- Germline murine immunoglobulin IGHV genes in wild-derived and classical inbred strains: a comparison

- Organization of the variable region of the immunoglobulin heavy-chain gene locus of the rat.

- IGHV1, IGHV5 and IGHV7 subgroup genes in the rhesus macaque.

1    **TITLE: A comparison of immunoglobulin IGHV, IGHD and IGHJ genes in wild-**

2    **derived and classical inbred mouse strains**

3

4    **Short running title: Germline IGH genes in inbred mice**

5

6    Corey T. Watson[1]*, Justin T. Kos[1], William S. Gibson[1], Leah Newman[2,3], Gintaras

7    Deikus[2,3], Christian E. Busse[4], Melissa Laird Smith[2,3], Katherine J.L. Jackson[5], Andrew

8    M. Collins[6]*

9

10   [1]Department of Biochemistry and Molecular Genetics, University of Louisville School of

11   Medicine, Louisville, KY USA 40202; [2]Icahn Institute for Genomics and Multiscale

12   Biology, Icahn School of Medicine at Mount Sinai, New York, NY 10029; [3]Department of

13   Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York,

14   NY 10029; [4]Division of B Cell Immunology, German Cancer Research Center, 69120

15   Heidelberg, Germany; [5]Immunology Division, Garvan Institute of Medical Research,

16   Darlinghurst, 2010 New South Wales, Australia; [6]School of Biotechnology and

17   Biomolecular Sciences, University of New South Wales, Sydney, 2052 New South

18   Wales, Australia

19   *To whom correspondence should be addressed.

20   Corey T. Watson: corey.watson@louisville.edu

21   Andrew Collins: a.collins@unsw.edu.au

24 **ABSTRACT**

25  The genomes of classical inbred mouse strains include genes derived from all three major

26  subspecies of the house mouse, *Mus musculus*. We recently posited that genetic diversity in the

27  immunoglobulin heavy chain (IGH) gene loci of C57BL/6 and BALB/c mice reflect differences in

28  subspecies origin. To investigate this hypothesis, we conducted high-throughput sequencing of

29  IGH gene rearrangements to document IGH variable (IGHV), joining (IGHJ), and diversity

30  (IGHD) genes in four inbred wild-derived mouse strains (CAST/EiJ, LEWES/EiJ, MSM/MsJ, and

31  PWD/PhJ), and a single disease model strain (NOD/ShiLtJ), collectively representing genetic

32  backgrounds of several major mouse subspecies. A total of 341 germline IGHV sequences were

33  inferred in the wild-derived strains, including 247 not curated in the International

34  Immunogenetics Information System. In contrast, 83/84 inferred NOD IGHV genes had

35  previously been observed in C57BL/6 mice. Variability among the strains examined was

36  observed for only a single IGHJ gene, involving a description of a novel allele. In contrast,

37  unexpected variation was found in the IGHD gene loci, with four previously unreported IGHD

38  gene sequences being documented. Very few IGHV sequences of C57BL/6 and BALB/c mice

39  were shared with strains representing major subspecies, suggesting that their IGH loci may be

40  complex mosaics of genes of disparate origins. This suggests a similar level of diversity is likely

41  present in the IGH loci of other classical inbred strains. This must now be documented if we are

42  to properly understand inter-strain variation in models of antibody-mediated disease.

43

44

45

46

47

48

49 **INTRODUCTION**

50 Inbred mouse strains are critical to biomedical research, and many of the most

51 important strains such as DBA, C57BL, C3H, CBA and BALB/c have now been in use

52 for almost a century [1]. The C57BL and BALB/c strains have been particularly important

53 for our understanding of the biochemistry and immunogenetics of immunoglobulins (IG)

54 [2-4]. This understanding was achieved despite a lack of detailed knowledge of the

55 antibody genes of these and other inbred laboratory mouse strains, and until recently it

56 was thought that the genes present in these different strains were likely to be highly

57 similar [5].

58

59 Mouse antibody genes were first identified using cell lines derived from BALB/c mice

60 because of the availability of mineral-oil induced plasmacytomas from this strain [4]. The

61 cataloguing of BALB/c antibody genes effectively ceased with the emergence of the

62 C57BL/6 strain as the workhorse of transgenic and genomic studies. The IG heavy

63 chain region (IGH) locus of the C57BL/6 strain was therefore the first to be sequenced

64 and annotated [6, 7]. Part of the IGH locus of the 129S1 strain was also reported [8], before

65 comprehensive genomic investigation of mouse germline IGHV genes essentially

66 ceased. By this time, two databases had catalogued mouse IGH genes and apparent

67 allelic variants of these genes: the VBASE2 [9] and IMGT [10] databases. A positional

68 nomenclature was then developed by IMGT, based upon the mouse genome reference

69 sequence, while an alternative positional nomenclature was developed by Johnston and

70 colleagues, based upon an alternative assembly of the C57BL/6 genome [6]. A non-

71 positional gene sequence identifier system was also developed by VBASE2 [9].

3

72

73     The study of IGH gene variation in humans followed a similar trajectory to that of the

74     mouse. Genes and likely allelic variants were reported over a twenty-year period,

75     starting in the late 1970s [11]. There was a sharp decline in the reporting of new

76     sequences once the complete human IGH locus was published in 1998 [12], but the

77     advent of high-throughput sequencing of human antibody genes reawakened interest in

78     the documentation of allelic variants [13, 14]. A surprising level of antibody gene variation,

79     including structural variation of the IGH locus, has since been shown within the human

80     population [15-18], and such variation can have important consequences for the

81     development of a suitable protective antibody repertoire [19, 20]. A similar exploration of the

82     immunoglobulin gene variation and antibody repertoires is now beginning in the mouse.

83

84     Many recently discovered allelic variants of IGH variable (IGHV), diversity (IGHD), and

85     joining (IGHJ) genes have been identified from sets of VDJ gene rearrangements, using

86     a process of inference [21]. Rearranged VDJ sequences are often affected by somatic

87     hypermutations, and such mutations are distributed throughout VDJ rearrangements.

88     When the same mismatch to a known germline gene is repeatedly observed in data

89     from a single subject, however, it is more likely that the nucleotide in question is a single

90     nucleotide variant (SNV), rather than being a nucleotide that has arisen by somatic

91     hypermutation [21]. When such mismatches are repeatedly seen in a large set of VDJ

92     rearrangements having diverse CDR3 regions and amplified from a single individual, the

93     inference of a previously undiscovered gene polymorphism may be made with

94     confidence. The discovery of allelic variants by inference is now a feature of many

4

95    human repertoire studies, and this is facilitated by a number of recently developed

96    utilities [22-25]. When this approach was applied to the mouse, the outcome was quite

97    unexpected.

98

99    Analysis of thousands of C57BL/6 VDJ rearrangements identified 99 of the 114

100   germline IGHV genes that have been reported to be functional in this strain [5]. It was

101   concluded that the remaining 16 genes are either non-functional or are expressed at

102   such low frequencies that they would only be detectable in deeper sequencing studies.

103   It was also concluded that all IGHV genes carried by any strain of inbred mouse that are

104   expressed at moderate frequencies should be readily determinable by inference from

105   VDJ gene datasets.

106

107   An analysis of BALB/c VDJ rearrangements was then performed, and 163 BALB/c IGHV

108   gene sequences were identified [5], only half of which were present in the IMGT database

109   [10]. The expression of ten unique IGHD sequences and four IGHJ genes was confirmed,

110   while three other reported BALB/c IGHD genes appeared to be non-functional. Although

111   the identification of BALB/c IGHV genes was almost certainly incomplete, the study

112   successfully captured the germline gene variability that accounts for almost all of the

113   genetic variation in the repertoire of rearranged heavy chain genes.

114

115   These were not the first striking differences seen in the IGH loci of BALB/c and C57BL/6

116   mice. Historically, these strains have been known to carry different IGH constant region

117   gene haplotypes, Igh-1[a] and Igh-1[b] respectively, and in the past these haplotypes were

118    determined serologically. Sequencing studies later showed striking differences between

119    the two haplotypes, particular at the IGHG2 gene locus. Controversy has surrounded

120    the evolutionary origins of this allotypic variation, but genomic evidence led to the

121    suggestion that it resulted from gene duplication and gene loss in different mouse

122    subspecies [26].

123

124    The differences, seen by Collins and colleagues [5], between the germline BALB/c and

125    C57BL/6 IGHV genes were so profound that it was proposed that the IGH genes of

126    these strains may have had their origins in different subspecies of the house mouse.

127    Three major subspecies of the house mouse have been described (*Mus musculus*

128    *musculus*, *M. m. domesticus* and *M. m. castaneus*) and subsequent analysis of genomic

129    SNV data supported the origins of BALB/c and C57BL/6 mice from *M. m. domesticus*

130    and *M. m. musculus* respectively [27] (see Table 1 and Supplementary figure 1).

131

132    To test the hypothesis that BALB/c and C57BL/6 IGH genes are derived from different

133    subspecies of the house mouse, in this study we first document IGHV, IGHD and IGHJ

134    germline genes in a number of wild-derived inbred mouse strains representing each of

135    the three major subspecies of the house mouse (CAST/EiJ: *M. m. castaneus*;

136    PWD/PhJ: *M. m. musculus*; LEWES/EiJ: *M. m. domesticus*), and of a wild-derived strain

137    (MSM/MsJ) that originated from *M. m. molossinus* mice that are generally considered to

138    be hybrids of *M. m. musculus* and *M. m. castaneus* (see Table I).  We then compare

139    these genes to those of the C57BL/6 and BALB/c strains, in order to infer the ancestry

140    of these classical inbred strains. The wild-derived strains used in this study were

6

141  developed in the 1970s from pairs of wild mice from known locations, with the intention

142  that each inbred strain would carry a genome that was derived from a single subspecies

143  of the house mouse. To explore the differences between strains that appear, based

144  upon preliminary SNV analysis, to be derived from the same subspecies of the house

145  mouse, we also investigated the NOD/ShiLtJ inbred strain (Table 1 and Supplementary

146  figure 1).

147

148  Although it was not possible in this study to document all rearrangeable genes of the

149  IGH loci, coverage was sufficient to allow broad conclusions to be drawn. Unexpectedly,

150  the NOD and C57BL/6 gene IGHV, IGHD, and IGHJ loci appear to be almost exactly

151  the same, while striking differences exist between each wild-derived strain and both the

152  C57BL/6 and BALB/c strains. The divergence that was seen between the strains

153  suggests that the IGH loci of inbred mouse strains are likely to harbor so much genetic

154  variation that if the antibody responses of these and other classical inbred mouse

155  strains are to be properly understood, it will be essential to fully document their germline

156  genes. The divergence between the strains also suggests that a single positional mouse

157  gene nomenclature based upon the C57BL/6 genome reference sequence may fail to

158  properly represent the genes of many important inbred mouse strains.

159

160  **RESULTS**

161  ***Defining inferred IGH germline gene sets in diverse inbred mouse strains***

162  To ensure the highest quality input data, prior to VDJ assignment we leveraged the

163  PacBio circular consensus (CCS2) algorithm to generate high quality circular consensus

7

164     reads for each sample. The average read length across the libraries sequenced was

165     22.5 Kb (Supplementary table 1). The long read lengths paired with the target amplicon

166     library size of ~1200 bp, resulted in a mean of 23.9 circular consensus passes per

167     amplicon (Supplementary table 1). Finally, we applied a Q30 cutoff to data from each

168     library, resulting in a total of 36782, 43522, 43173, 28136 and 43044 pre-mapped raw

169     reads for CAST/EiJ, LEWES/EiJ, MSM/MsJ, PWD/PhJ, and NOD/ShiLtJ, respectively,

170     with mean CCS read scores of 1 (Supplementary table 1). Each of these high quality

171     read datasets were used for IGHV, IGHD, and IGHJ gene assignment, clonal

172     assignment, and germline gene inference. Read data for each strain have been

173     submitted to the Sequence Read Archive (SRA) under the BioProject ID PRJNA533312.

174

175     The discovery stage of our inference pipeline for each strain was modelled after that

176     used previously [5] (see Materials and Methods below). Clonal assignment and clustering

177     resulted in a total of 5261 (CAST/EiJ), 5042 (LEWES/EiJ), 4374 (MSM/MsJ), 3827

178     (PWD/PhJ), and 3743 (NOD/ShiLtJ) unique clones; a single representative sequence

179     from each of the identified clones was then randomly selected to use for germline

180     inference. With these sequences as input, a total of 87, 78, 84, 92, and 84 germline

181     IGHV sequences were inferred for CAST/EiJ, LEWES/EiJ, and MSM/MsJ, PWD/PhJ,

182     and NOD/ShiLtJ, respectively (Figure 1a; Supplementary table 2). Each inferred

183     germline sequence was represented by at least 0.1% of the total clones observed in a

184     given strain; the numbers of clones representing each inference are provided in

185     Supplementary table 2 and plotted in Supplementary figure 2.

186

187    Across the wild-derived strains, the validity of some inferred germline IGHV sequences

188    was supported by their presence in the IMGT mouse reference directory

189    (http://www.imgt.org) (Figure 1a). This included 5 sequences in CAST/EiJ, 15 in

190    LEWES/EiJ, 62 in MSM/MsJ, and 12 in the PWD/PhJ strain. Sequences inferred from

191    the wild-derived strains were, however, dominated by non-IMGT alleles (247/341, 72%;

192    Figure 1a). Of special note, all but one of the NOD/ShiLtJ gene sequences have

193    previously been reported in the IMGT reference directory, and all of these 83 sequences

194    are reported there as C57BL/6 sequences.

195

196    Additional supporting evidence was found in other public sequence repositories for 44

197    non-IMGT sequences (Supplementary table 2). Two non-IMGT sequences were found

198    in the MSM/MsJ strain with perfect matches to sequences reported in either VBASE2

199    (http://www.vbase2.org/) or the NCBI reference set. This was also true for 2, 29, and 11

200    non-IMGT sequences inferred from CAST/EiJ, LEWES/EiJ, and PWD/PhJ, respectively.

201

202    The sets of germline inferences were used as starting databases for analysis by

203    IgDiscover [22], to seek further validate the inferences. This analysis confirmed 93% of the

204    inferences: 82/87 (94%), 68/78 (87%), 76/84 (90%), 86/92 (93%), and 84/84 (100%) of

205    the inferred sequences for CAST/EiJ, LEWES/EiJ, and MSM/MsJ, PWD/PhJ, and

206    NOD/ShiLtJ, respectively.

207

208    Most of the inferences in the wild-derived strains that were not confirmed by IgDiscover

209    were present at low copy number (Supplementary table 2 and Supplementary figure 3),

9

210    though a few unconfirmed sequences were seen at relatively high copy number,

211    including CAST-IGHV9-2 (28 clones), LEWES-IGHV3-2 (25 clones) and MSM-IGHV2-4

212    (50 clones). Many of the unconfirmed sequences were supported by other secondary

213    evidence from public databases, or were also seen in other strains from this study

214    (Supplementary table 2). This included seven sequences that are present among

215    whole-genome shotgun sequence data generated by the Mouse Genome Project

216    (https://www.sanger.ac.uk/science/data/mouse-genomes-project). Just 11 of the

217    sequences unconfirmed by IgDiscover lacked additional evidence supporting their

218    validity, representing only ~2.5% of the total number of sequences identified in our

219    dataset. Issues of false negatives by IgDiscover, and other related inference tools, have

220    been reported, but have not been fully explained [24]. As a consequence, after additional

221    manual review of our results, the full sets of inferences were retained and are reported

222    here.

223

224    In some cases, novel sequences from wild-derived strains were quite divergent from

225    published IGHV sequences, varying from 87.02% to 99.66% sequence identity. This

226    was also dependent on strain in that, among the wild-derived strains, inferred germlines

227    from CAST/EiJ, LEWES/EiJ, and PWD/PhJ exhibited much greater sequence

228    divergence from IMGT alleles than sequences inferred in MSM/MsJ (Figure 1b).

229

230    The counts of genes within the different IGHV families were generally comparable

231    across the five strains (Figure 1c), with some exceptions. The germline repertoires of all

232    strains were clearly dominated by the IGHV1 family. However, the numbers of IGHV

233    genes in other subgroups were more variable. For example, the repertoire of CAST/EiJ

234    harbored fewer IGHV2 genes relative to the other strains, but greater numbers of

235    IGHV3, IGHV6, IGHV9, and IGHV14 genes. At least one representative germline

236    sequence of subfamilies IGHV1-IGHV3, IGHV5-IGHV10, and IGHV14 was inferred from

237    all strains; in contrast, sequences of the remaining subfamilies, which are all small

238    subfamilies in the C57BL/6 strain, were absent in at least one strain. IGHV13 and

239    IGHV15 sequences were only observed in two of the five strains. Whether this

240    represents a genuine lack of functional IGHV subfamily sequences in these strains (e.g.

241    as a result of pseudogenization or genomic deletion), or whether this is due to under

242    sampling of these repertoires is not clear. Deeper sequencing and ultimately

243    comprehensive genomic characterization in each strain would be needed to fully assess

244    this. Consistent with the general subfamily distributions observed, the majority of non-

245    IMGT sequences in CAST/EiJ (n=32), LEWES/EiJ (n=32), and PWD/PhJ (n=39) were

246    represented by IGHV1 genes (Figure 1d). In contrast, however, although the MSM/MsJ

247    repertoire was also dominated by IGHV1, the majority of non-IMGT sequences

248    identified in that strain were from IGHV2 (n=9) and IGHV5 (n=9) (Figure 1d).

249

250    Analysis was also performed to identify previously unreported polymorphism of IGHD

251    and IGHJ genes, and to define the sets of IGHD and IGHJ genes in each mouse strain.

252    Analysis of IGHD gene alignments led to the identification of eight or nine IGHD genes

253    in each of the strains, with four novel sequences being identified as likely allelic variants

254    of IGHD1-1 (n=2), IGHD2-3 (n=1) and IGHD2-12 (n=1) (Figure 2a and Supplementary

255    table 3). These polymorphisms could be determined with absolute certainty because the

11

256    critical nucleotides that distinguish them from previously reported IGHD sequences are

257    sufficiently distant from the IGHD gene ends. No variants were identified with

258    differences in the 5' or 3' terminal nucleotides of the IGHD genes, but we cannot

259    exclude the possibility that our analysis failed to detect such variation.

260

261    The IGHD genes of NOD/ShiLtJ were shown to be identical to those of the C57BL/6

262    mouse. The MSM/MsJ and PWD/PhJ strains share 6 IGHD genes with the C57BL/6

263    strain, and carry two and three additional unique genes respectively. The PWD/PhJ

264    strain expresses a variant of the 26 nucleotide IGHD2-12 gene (Figure 2a and

265    Supplementary table 3). The LEWES/EiJ strain shares five sequences with the BALB/c

266    strain, but also expresses four additional sequences, including a previously unreported

267    variant of the IGHD1-1 gene. This variant differs from the IGHD1-1*01 sequence at two

268    positions (Figure 2a and Supplementary table 3). The CAST/EiJ strain also expresses a

269    novel IGHD1-1 variant, distinct from that observed in LEWES/EiJ, as well as six genes

270    that are expressed by other strains in the study. Only CAST/EiJ expresses the IGHD2-

271    13*01 gene. The IGHD4-1*01 and IGHD4-1*02 are highly similar gene segments

272    sharing 9 of 11 and 9 of 10 nucleotides, respectively. In all strains, full length examples

273    of both IGHD4-1 variants are observed, but the shorter IGHD4-1*02 is difficult to confirm

274    with confidence as such a sequence can be derived by trimming 2 nucleotides from the

275    5' end of  *01, followed by the non-template addition of 'c'.

276

277    All strains were found to carry four functional IGHJ genes, which were each used at

278    sufficient frequency to make their identification unequivocal (Figure 2b and

12

279 Supplementary table 4). Only in CAST/EiJ was there evidence of a novel allele

280 (IGHJ2_var). This novel allele harbored a single nucleotide difference from the closest

281 known allele in the IMGT reference directory (Figure 2b). The other CAST/EiJ IGHJ

282 genes were shared with the other strains and were identical to the genes carried by

283 C57BL/6 mice. No strain carried the BALB/c IGHJ1*01 allele [5].

284

285 ***Extensive IGHV germline diversity and limited overlap between strains***

286 We next investigated the extent of overlap of IGHV sequences among the surveyed

287 strains. The sets of germline genes inferred from each inbred strain were compared to

288 sequences identified in all other strains to determine how many sequences were

289 identical between strains. Comparisons were additionally made with previously

290 published inferences from BALB/c mice (n=163) [5], and with the IMGT repertoire of

291 functional C57BL/6 sequences (n=114). Surprisingly little overlap was observed

292 between strains, and the majority of inferred germline sequences were unique to a

293 single strain (Figure 3a). Among the wild-derived strains surveyed, CAST/EiJ had the

294 highest number of unique germline sequences (76/87 sequences). On the other hand,

295 83 shared sequences were observed in NOD/ShiLtJ and C57BL/6, with 48 of these 83

296 sequences being additionally shared by MSM/MsJ. A single sequence was identified in

297 five different strains (LEWES/EiJ, NOD/ShiLtJ, CAST/EiJ, PWD/PhJ, and C57BL/6;

298 Supplementary table 2).

299

300 We further explored interstrain IGHV sequence relationships by estimating the average

301 sequence similarities of IGHV sets between strains. Consistent with sequence overlaps

302  presented in Figure 3a, we noted a range of mean pairwise sequence identities,

303  depending on the strains in question. For example, sequences in MSM/MsJ, C57BL/6,

304  and NOD/ShiLtJ, strains which share the most identical sequences with one another

305  (Figure 3a), also have high average pairwise sequence identities (>99%; Figure 3b; see

306  also Supplementary figure 4 for full pairwise comparisons). This is in contrast to mean

307  identities observed for all other pairwise strain comparisons, which ranged from 94.8%

308  to 97.1%. These levels of identity also generally matched the number of shared

309  sequences between strains. For example, LEWES/EiJ shared the most identical

310  sequences with BALB/c, and among all pairwise sequence comparisons between

311  LEWES/EiJ and other strains, the highest mean sequence identity was with the BALB/c

312  germline set (97.1%; Figure 3b).

313

314  No attempt was made in this study to determine whether or not any pairs of highly

315  similar sequences could be allelic variants of a single gene. In light of known structural

316  variation in the IGH loci of C57BL/6 and BALB/c mice [7], we must assume that significant

317  structural variation is possible between the loci of the wild-derived strains reported here.

318  In a locus where gene duplications and deletions, as well as pseudogenization, are

319  central drivers of evolution, and where many sets of highly similar genes are found in

320  each strain, the documentation of allelic variants of any gene will require comprehensive

321  genomic sequencing.

322

323  ***Implications of inferred germline references on repertoire alignments***

14

324   The re-analysis of the VDJ datasets from each strain using the strain-specific inferred

325   repertoires dramatically impacted upon the repertoire-level alignment metrics (Table 2).

326   This was especially true for CAST/EiJ which was most poorly represented in the IMGT

327   reference directory. Considering the IgBLAST output for the two different repertoires,

328   just 3.8% of IGHVs and 49.4% of IGHJs were unmutated if the IMGT reference directory

329   was used for analysis. This increased to 87.8% and 88.2% for IGHV and IGHJ

330   respectively using the CAST/EiJ inferred germline references. LEWES/EiJ and

331   PWD/PhJ also experienced significant improvements for the IGHV, shifting the

332   frequency of unmutated reads within the IgM repertoires from 35.6 to 88.6% and 13.9 to

333   82.9%, respectively.

334

335   **DISCUSSION**

336   This study was undertaken to investigate the hypothesis that IGH genes of subspecies

337   of the house mouse are highly divergent, and to trace the subspecies origins of the IGH

338   loci of BALB/c and C57BL/6 inbred mice. Since the mice that were used to establish the

339   classical inbred strains of laboratory mice came from diverse and usually

340   undocumented sources, we reasoned that differences in the loci of the various

341   subspecies of the house mouse could explain the marked differences that were

342   previously reported between the sets of IGHV genes found in the BALB/c and C57BL/6

343   strains. To test this hypothesis, we inferred the germline IGHV, IGHD and IGHJ genes

344   of wild-derived strains representing each of the three major subspecies of the house

345   mouse (CAST/EiJ: *M. m. castaneus*; PWD/PhJ: *M. m. musculus*; LEWES/EiJ: *M. m.*

346   *domesticus*), and of a wild-derived strain (MSM/MsJ) that originated from *M. m.*

15

347    *molossinus* mice that are generally considered to be natural hybrids of *M. m. musculus*

348    and *M. m. castaneus* (see Table I).  Whereas SNV-inferred haplotypes reported by

349    Yang and colleagues [28] supported the reported subspecies origins of CAST/EiJ,

350    PWD/PhJ, and LEWES/EiJ, these data suggested that the MSM/MsJ IGH locus is *M. m.*

351    *musculus*-derived, with a SNV profile that is little different to that of the C57BL/6 mouse

352    (see Table I and Supplementary figure 1). The IGH locus of the MSM/MsJ strain was

353    therefore investigated in the hope that it would shed light on genetic variation within the

354    loci of *M. m. musculus*-derived strains. For the same reason, the NOD/ShiLtJ strain was

355    also included in this study, because SNV analysis suggested that it too carries a

356    C57BL/6-like IGH haplotype.

357

358    In general, among the wild-derived strains surveyed in this study we observed

359    surprisingly little overlap between IGHV sequences. While this was expected in

360    comparisons of strains predicted to carry IGH loci originating from different subspecies,

361    intriguingly there were notable differences in IGHV sequence sets between strains

362    carrying loci of the same predicted subspecies. For example, PWD/PhJ mice only

363    shared 13 (~14%) IGHV gene sequences with the C57BL/6, NOD/ShiLtJ, and MSM/MsJ

364    strains, despite the fact that all of these strains were predicted to share IGH loci of *M.*

365    *m. musculus* origin. In fact, the PWD/PhJ strain shared almost the same number of

366    IGHV sequences with the CAST/EiJ, LEWES/EiJ, and BALB/c strains (predicted to

367    represent the *M. m. castaneus, M. m. domesticus* and *M. m. domesticus* subspecies

368    respectively). PWD/PhJ IGHV sequences were also no more similar to sequences of

369    other *M. m. musculus*-derived strains than to *M. m. castaneus*- or *M. m. domesticus*-

16

370 derived strains. Similarly, although SNV analysis suggested that BALB/c and

371 LEWES/EiJ mice both carry *M. m. domesticus*-derived IGH loci, only 13 (~16%) of the

372 identified LEWES/EiJ IGHV sequences are shared with the BALB/c strain. LEWES/EiJ

373 IGHV sequences were collectively most similar to BALB/c genes, relative to the other

374 strains sequenced here, but it is difficult to believe that the IGH loci of both strains are

375 derived in their entirety from shared *M. m. domesticus* ancestors.

376

377 The IGHV gene sequences of the MSM/MsJ strain were particularly surprising. Few of

378 the 84 MSM sequences were seen in any of the other three wild-derived strains. As

379 expected, none of these sequences were amongst the 78 IGHV genes identified in the

380 *M. m. domesticus*-derived LEWES/EiJ strain; however, there was also little identity with

381 sequences identified in either of the subspecies that are said to have given rise to the

382 hybrid *M. m. molossinus* mice. Only 10 of the 84 MSM/MsJ IGHV sequences matched

383 those seen in the *M. m. musculus*-derived PWD/PhJ strain, and only 4 sequences

384 matched those in the *M. m. castaneus*-derived CAST/EiJ strain. Instead, substantial

385 identity was seen between MSM/MsJ mice and inbred C57BL/6 and NOD/ShiLtJ mice,

386 that SNV analysis suggests are both *M. m. musculus*-derived (Supplementary figure 1).

387 Genomic sequencing will be required to determine whether or not this identity is a

388 consequence of an unreported breeding accident in the history of the MSM/MsJ colony.

389

390 Analysis of the expression of IGHJ genes showed little variation between strains,

391 though interestingly, no strain shared the expression of the BALB/c IGHJ1*01

392 sequence. Analysis of IGHD gene expression was more informative. The NOD/ShiLtJ

17

393    mice appear to carry an IGHD locus that is identical to that of the C57BL/6 strain. Not

394    only does this analysis support the conclusions from the IGHV gene analysis, but it

395    gives confidence in the inference process itself. The IGHD genes of the MSM strain

396    were also similar to those of the C57BL/6 strain, with seven of the nine C57BL/6 IGHD

397    sequences being shared. Against expectations, just five of the ten BALB/c IGHD

398    sequences were shared by *M. m. domesticus*-derived LEWES mice, while 8 sequences

399    were shared with the *M. m. musculus*-derived PWD strain. We cannot rule out the

400    possibility that this analysis failed to identify SNVs in either the terminal 5' or 3'

401    nucleotides of the IGHD sequences, though we believe this is unlikely. It has been

402    shown that the inference process is often unable to accurately identify terminal

403    nucleotides in AIRR-Seq data [29]. On the other hand, the antibody repertoire of the

404    mouse is strongly shaped by joining at sequences of short homology [30] and there may

405    be strong evolutionary pressure maintaining, for example, the identity between the 3'

406    ends of the IGHD2 gene family and the 5' ends of some IGHJ sequences, thereby

407    preventing the development of variation in the terminal IGHD gene nucleotides.

408

409    The demonstration that NOD/ShiLtJ mice carry IGHV, IGHD and IGHJ loci that are so

410    closely related to those of the C57BL/6 strain is intriguing and also quite unexpected,

411    because no direct relationship between these strains has been reported. The

412    NOD/ShiLtJ mouse was derived in the 1970s from cataract-prone CTS mice, which in

413    turn were developed in the 1960s from outbred Swiss mice [31, 32]. C57BL/6 mice, on the

414    other hand, were developed in the 1920s by Clarence Little from the progeny of a pair of

415    'fancy mice' [33, 34]. It is difficult to believe that the NOD and C57BL/6 loci could have

416   arisen independently by the chance selection of unrelated outbred founder pairs. The

417   nearly identical sets of IGHV genes in the C57BL/6 and NOD/ShiLtJ mice are more

418   suggestive of introgression, with a C57BL/6-like locus being introduced into the

419   ancestors of the modern NOD strain by outcrossing. This notion is further supported by

420   the observation that NOD clusters together with C57BL/6, based on mitochondrial SNV,

421   but is rather distant from the strain based on chromosome Y SNV [28]. This hints of a

422   contamination via the maternal line. Of note, this would not be the first reported

423   breeding accident involving the NOD lineage [35].

424

425   If the IGHV and IGHD genes of the LEWES/EiJ, PWD/PhJ and CAST/EiJ mice are

426   accepted as being broadly representative of the three major subspecies of the house

427   mouse, then neither the BALB/c strain nor the C57BL/6 strain can be unequivocally

428   linked with one or other of the three major subspecies of the house mouse. It may be

429   that the IGH loci of these and other classical inbred strains have a mosaic structure,

430   representing alternating blocks of genes that are potentially polymorphic within one of

431   the three major mouse subspecies, or divergent between them. It is also possible that

432   the IGH loci of these strains include haplotype blocks derived from other lineages of the

433   house mouse, or even from other Mus species such as *M. spretus*. Other subspecies of

434   *M. m. musculus* probably exist, and a number have been proposed, such as *M. m.*

435   *bactrianus* and *M. m. gentilulus* [36, 37].

436

437   The divergence patterns of the IGH loci of the mouse strains reported here can be

438   considered in the context of the allelic diversity that has been reported in humans. For

439   example, a comparison of the two published fully sequenced human IGHV haplotypes

440    revealed that, of the 68 non-redundant functional/ORF IGHV sequences identified

441    across the two haplotypes, only 22 were observed in both, accounting for both allelic

442    and structural variants [18]. In addition, although the population genetics of the human

443    IGH locus remains relatively uncertain [18-20], some indirect measures of diversity are

444    available to us. One measure of diversity is provided by consideration of heterozygous

445    loci in individuals who have been genotyped by the analysis of VDJ rearrangements.

446    When heterozygosity was explored at 50 IGHV gene loci in 98 individuals, only 5 genes

447    were heterozygous in more than 50% of individuals, and only 19 genes were

448    heterozygous in more than 20% of individuals [15]. However, it is notable that more

449    extreme examples of IGHV heterozygosity have been reported in some populations; for

450    example, a study of genomic DNA in 28 South Africans characterized >50 IGHV alleles

451    in all but two individuals [16]. This same study characterized 123 alleles that were not

452    present in IMGT. Similarly, a study of 10 individuals from Papua New Guinea identified

453    17 previously unreported IGHV sequences; however, in each individual, all but two or

454    three sequences had previously been reported from studies in Europe, America and

455    Australia [17]. Collectively, these data suggest that there are many alleles common across

456    different populations, but in some cases, it is likely that some alleles will also be more or

457    less frequent in, or even private to specific populations. Taken together, it seems likely

458    that the level of IGH diversity in inbred laboratory mice could be more extensive than

459    what has been observed in human studies conducted to date. This may be expected

460    given that our data suggest that IG genes present in many of the strains in use today

461    are likely to harbor variation originating from many different subspecies.

462

463    If we are to properly understand the antibody repertoires of the major laboratory strains,

464    their IGH loci will all need to be separately investigated. Unfortunately, existing projects

465    such as the Mouse Genome Project (MGP) are unlikely to provide the necessary data.

466    The MGP is currently sequencing the genomes of 16 inbred mouse strains, but it is

467    unlikely that such assemblies will lead to suitably reliable reference sequences for the

468    IGH locus in the short term, particularly without targeted and focused efforts that

469    attempt to deal with the complexity and associated technical issues of these regions.

470    Thus, immediate advances in our knowledge of the IGH loci of inbred strains will likely

471    come from the inference of genes from AIRR-seq data. The resulting lack of knowledge

472    of non-coding regions, and the lack of positional data, may make it difficult to decipher

473    relationships between some strains, but it should provide the basic information

474    regarding germline genes that is needed for accurate repertoire studies. Indeed, from

475    our reanalyses (Table 2) of data based on the novel germline sets reported here, it is

476    clear that strain-specific data will be needed in many cases if we are to ensure precision

477    in AIRR-seq studies of strains other than C57BL/6.

478

479    Deeper sequencing will also be required if a more complete documentation of the

480    available repertoire of germline genes in any inbred strain is to be determined.

481    Nevertheless, the partial documentation of the IGH loci reported here, and the strain

482    differences seen are sufficient to raise the possibility that the IGHV locus may contribute

483    to strain-related differences in mouse models of human disease. Allelic variants have

484    been associated with differences in disease susceptibility of rats [38] and humans [19]. IGHV

485    sequence variability might also contribute to the differences that have been reported in

486 the susceptibility of inbred mouse strains to both infectious [39, 40] and autoimmune

487 diseases [41].

488

489 The discovery of striking differences in the number of apparently functional IGHV genes

490 between the C57BL/6 and the BALB/c strains first raised the possibility that the

491 continued use of a positional nomenclature system for IGHV genes could be

492 problematic [5]. The results presented here confirm that this is the case. The C57BL/6

493 locus is unable to serve as a map for IGHV genes from other strains, and this appears

494 to be the case even for inbred strains that were shown in earlier SNV analysis to carry

495 *M. m. musculus*-derived IGH loci. A non-positional nomenclature should therefore be

496 developed. Attention should also be paid to the light chain loci carried by different

497 mouse strains. We have shown by SNV analysis that the three major subspecies of the

498 house mouse may all have contributed to the kappa loci of the major strains of inbred

499 laboratory mice [42]. If the kappa and lambda loci of laboratory mice are shown to have

500 the same kind of strain to strain variability as we have shown here for the IGH locus,

501 then a new non-positional nomenclature will be required for all the genes of the IG loci.

502

503 **METHODS**

504 ***Antibody gene repertoire sequencing***

505 Whole dissected spleens, preserved in RNAlater, were obtained from female mice from

506 Jackson Laboratories (https://www.jax.org) for five inbred strains (CAST/EiJ [JAX stock

507 #000928], n=1; LEWES/EiJ [JAX stock #002798], n=1; PWD/PhJ [JAX stock #004660],

508 n=1; MSM/MsJ [JAX stock #003719], n=1; NOD/ShiLtJ [JAX stock #001976], n=1). An

509  individual mouse was studied from each strain, as the power of this kind of investigation

510  comes from sequencing depth rather than from the investigation of biological replicates.

511  Each VDJ rearrangement provides independent support for the presence of a gene

512  segment in the genome.

513

514  Total RNA was extracted from a section of each spleen using the RNeasy Mini kit

515  (Qiagen, Cat. No. 74104; Germantown, MD). For each strain, 5' RACE first-strand

516  cDNA synthesis was conducted using the SMARTer RACE cDNA Amplification Kit

517  (Takara Bio, Cat. No. 634858; Mountain View, CA), with an input amount of 1 µg of

518  RNA per sample. Rearranged VDJ-IgM amplicons were generated using a single IgM

519  oligo positioned in the CH3 region of the mouse IGHM gene (5'-

520  CAGATCCCTGTGAGTCACAGTACAC-3'; 10 µM), paired with a universal primer

521  (Takara Bio; 5'-AAGCAGTGGTATCAACGCAGAGT-3'; 10 µM). First-strand cDNA for

522  each strain was amplified using Thermo Fisher Phusion HF Buffer (Thermo Fisher, Cat.

523  No. F530S; Waltham, MA) for 30 PCR cycles. Amplicons were run on 2% agarose gel

524  for size confirmation. Final amplicons were used to generate SMRTbell template

525  libraries and each library was sequenced across 2 SMRTcells on a Pacific Biosciences

526  RSII system using P6/C4 chemistry and 240 minute movies (Pacific Biosciences; Menlo

527  Park, CA).

528

529  **Data processing and germline gene inference**

530  Reads from each RSII run were combined and processed using the CCS2 algorithm.

531  CCS2 reads of Q30 or greater were processed with pRESTO [43] v0.5.1 as follows: 1)

23

532    Reads without a universal primer were removed using a maximum primer match error

533    rate of 0.2 and a maximum template-switch error rate of 0.5 to de-duplicate the reads;

534    2) Reads without an IgM primer were removed using a maximum primer match error

535    rate of 0.2 and a maximum template-switch error rate of 0.5; 3) Duplicate sequences

536    from the FASTQ files were removed using the default value of a maximum of 20

537    ambiguous nucleotides. Following pRESTO processing, IGHV gene assignments were

538    noted after mapping reads to the IMGT germline database using IMGT/HighV-QUEST

539    version 1.6.0 (20 June 2018) [44]. Resulting IMGT summary output was processed using

540    Change-O v0.4.3 [45]. Reads representing incomplete or non-productive VDJ

541    rearrangements were discarded. All sequences in each sample were assigned to clones

542    using the distToNearest function in the SHazaM R package and the DefineClones

543    function in Change-O [45]. A single member of each clonal group was randomly selected

544    for downstream analyses.

545

546    For germline IGHV gene inference for each strain, the sequences representing the

547    strain's clonal groups were first clustered based on the IGHV gene assignment and

548    associated percent identity of the alignment to the nearest IMGT reference directory

549    gene sequence. Sequences shorter than the 138 base pair (bp) length of the shortest

550    reported functional IGHV sequence in IMGT were excluded from the analysis.

551    Consensus IGHV gene sequences were then determined for each cluster using CD-HIT

552    (cd-hit-est v4.6.8) [46] requiring that a given inference be represented by at least 0.1% of

553    total clonal groups identified in the strain's dataset. To be conservative, consensus

554    sequences for each inferred germline were trimmed to the shortest sequence

24

555   representative in the identified cluster. If the percent match identity to the closest IMGT

556   germline gene sequence was <95%, then the alignments were manually inspected for

557   evidence of chimeric PCR amplification [47].

558

559   An inferred germline IGHV sequence reference dataset was generated for each strain.

560   Predicted germline IGHV gene databases from Q30 datasets were assessed using one

561   iteration of IgDiscover v0.10 [22]. For the analysis of each strain, inferred germline IGHV

562   sequences generated from our clustering method were used as the IGHV gene

563   database for IgDiscover. Changes made to the default configuration were as follows: 1)

564   'race_g' set to 'true' to account for a run of G nucleotides at the beginning of the

565   sequence 2) 'stranded' set to 'false' because the forward primer was not always located

566   at the 5' end, and 3) 'ignore_j' set to 'true' to ignore whether or not a joining (J) gene

567   had been assigned for a newly discovered IGHV gene.

568

569   All      inferred      germline      sequences      were      compared,      using      BLAST

570   (https://blast.ncbi.nlm.nih.gov/) [48], to sequences in the following databases: IMGT

571   (http://www.imgt.org), VBASE2 (http://www.vbase2.org/), and the NCBI non-redundant

572   nucleotide sequence collection. Germline sets were compared across strains using

573   BLAT [49]. Perfect matches between sequences of different strains required full length

574   alignments of the query sequence at 100% identity; sequence length variation at the 3'

575   end of query and subject sequences was allowed. Calculations of mean IGHV

576   sequence identities between strains were computed by taking the mean of the

577   sequence identities for the best pair-wise hits of all genes in the smaller germline set of

578   two strains being compared. For the comparison between IGHV sequences of

25

579    NOD/ShiLtJ and MSM/MsJ, which had an equal number of inferred germline genes, we

580    opted to present the mean sequence identity using the NOD/ShiLtJ germline set.

581

582    Germline IGHD genes were inferred for each strain by further analysis of the junction

583    sequences from each strain's clone group representatives. Junctional sequences as

584    defined by IMGT/HighV-QUEST from the Change-O tables were compared against the

585    full IMGT mouse IGHD reference directory (Release 201918-4; 9 May 2019). Reference

586    directory genes that share identical coding sequences were grouped to a single entry

587    for comparison to the VDJ junctions. For each reference IGHD gene all substrings of

588    length 5 or more were compared to the all junction sub-sequences of the same length.

589    Matches between the IGHD and junctional substrings were retained if the hamming

590    distance was 2 or less, with the exception of rejecting the alignment if the pattern of

591    mismatches fell into either of the following categories; the mismatches were in the first

592    or final two positions of the substrings, or the third and fourth positions were both

593    mismatched, or the third and fourth to last positions were both mismatched. These

594    heuristics were implemented to avoid patterns of mismatch that arise from mis-

595    alignment of non-template encoded N additions in circumstances where the IGHD has

596    undergone exonuclease trimming. For each VDJ junction the set of IGHD substrings

597    that maximised length and minimising mismatches were included in the analysis.

598    Multiple IGHD substring matches were permitted per VDJ junction.

599    Counts for observed IGHD substring matches across each strain's dataset of clone

600    representations were tallied for analysis. Analysis considered all IGHD substring

601    matches with fewer than 10 exonuclease removals and with minimum lengths of 10

26

602     nucleotides (or 9 for IGHD4-4). The presence of a reference IGHD was accepted if

603     abundant, full-length matches were observed within a strain's datasets (generally in

604     excess of 10 observations). For reference gene-derived substrings that included

605     mismatches, if the same mismatch was observed across 80% of the substrings for the

606     IGHD gene then this was taken as putative evidence of the presence of SNVs in the

607     germline gene. If the mismatches could be explained by mis-assignment of another

608     IGHD which was already confirmed by full-length, perfect matches then the putative

609     variant was dismissed. In the absence of such mis-assignment and with the support of

610     the repeated detection of the SNV(s) associated with varied patterns of exonuclease

611     removal, a new IGHD variant was defined. For strains with new IGHD variants, the

612     substring analysis was repeated using the IMGT reference directory amended to

613     replace the IMGT gene(s) with the novel variant(s). Final tabulations were made from

614     these strain adjusted references.

615     IGHJ genes for each strain were inferred by analysis of the subset of clone

616     representatives for each strain that had unmutated IGHV regions in order to reduce the

617     likelihood of the IGHJ containing somatic point mutations. IGHJs that were more than 2

618     nucleotides shorter than the closest germline reference sequence were also excluded.

619     For each strain's dataset, IGHJ gene sequences were then clustered using CD-HIT,

620     requiring exact matches. Consensus sequences from each cluster were then examined

621     based on the number of supporting sequences. Within each strain, the top four clusters

622     based on their observed frequency in the repertoire were taken to represent the IGHJ

623     genes of that strain. Any additional clusters that had a frequency representing >10% of

624     one of the top four IGHJ gene clusters were also further examined. In only one case did

27

625     this occur (in CAST/EiJ); this sequence discarded because it was was determined to be

626     a derivative of one of the top four clusters from this strain, and to include divergent

627     bases originating from 5' NP addition.

628

629     Finally, each strain's dataset was re-aligned using strain-specific IGHV, IGHD and IGHJ

630     references. Strain specific references included the sequences listed in Supplemental

631     tables 2, 3 and 4. IMGT/HighV-QUEST was unable to be utilised for this analysis as it

632     does not permit specification of the non-IMGT germline reference sets. Alignments were

633     therefore performed with stand-alone IgBLAST [50] (version 1.14.0) with a reward score of

634     +1 and a mismatch penalty of -3. Scoring was adjusted in this way due to length

635     differences (283 - 303) for the inferred IGHVs due to conservative 3` end trimming. The

636     default +1/-1 scoring favours longer mismatched alignments, over shorter unmutated

637     alignments, even when the shorter alignments are full length with respect to the inferred

638     gene. The IgBLAST mouse auxiliary data file was amended to include the novel IGHJ2

639     variant from CAST/EiJ using the same offsets as IGHJ2*01. IgBLAST with the same

640     parameters was also used to align that VDJ datasets against the IMGT reference

641     directory to distinguish the contributions of the alignment algorithm and the germline

642     reference.

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659    **Table 1**: The intended *Mus musculus* subspecies origins of inbred strains of mice, and

660    the subspecies identity of the IGH loci of inbred strains, as determined by SNV analysis.

| Strain | Subspecies Origin | Subspecies identity by SNV analysis |
|---|---|---|
| PWD/PhJ | *M. m. musculus* | *M. m. musculus* |
| LEWES/EiJ | *M. m. domesticus* | *M. m. domesticus* |
| CAST/EiJ | *M. m. castaneus* | *M. m. castaneus* |
| MSM/MsJ | *M. m. mollosinus* | *M. m. musculus* |
| C57BL/6 | N/A | *M. m. musculus* |
| BALB/c | N/A | *M. m. domesticus* |
| NOD/ShiLtJ | N/A | *M. m. musculus* |

661

662

663

29

664

665

666

667

668

669

670

671

672

673

674 **Table 2:** IgBLAST results using the IMGT database compared to novel strain-specific

675 repertoires characterized in the current study.

676

| | | | IGHV | | | | IGHD | | | | IGHJ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | IGHV (n) * | unmutated IGHVs (n) ^ | unmuted reads (%) # | mean SHM (%) | IGHD (n) * | unmutated IGHDs (n) ^ | unmuted reads (%) # | mean SHM (%) | IGHJ (n) * | unmutated IGHJs (n) ^ | unmuted reads (%) # | mean SHM (%) |
| IgBLAST | IMGT | CAST | 109 | 30 | 3.8% | 2.21% | 16 | 16 | 87.9% | 0.92% | 8 | 4 | 49.4% | 1.45% |
| | | LEWES | 97 | 27 | 35.6% | 1.59% | 16 | 16 | 90.4% | 1.17% | 7 | 7 | 88.8% | 0.52% |
| | | MSM | 110 | 86 | 71.3% | 0.52% | 17 | 17 | 98.3% | 0.12% | 7 | 6 | 90.7% | 0.39% |
| | | NOD | 98 | 96 | 82.7% | 0.10% | 16 | 15 | 96.4% | 0.22% | 6 | 5 | 93.2% | 0.18% |
| | | PWD | 121 | 33 | 13.9% | 1.91% | 17 | 17 | 96.8% | 0.16% | 6 | 5 | 89.3% | 0.45% |
| IgBLAST | inferred | CAST | 87 | 87 | 87.8% | 0.25% | 9 | 9 | 96.2% | 0.30% | 4 | 4 | 88.2% | 0.49% |
| | | LEWES | 78 | 78 | 88.6% | 0.20% | 11 | 11 | 97.8% | 0.18% | 4 | 4 | 88.3% | 0.53% |
| | | MSM | 84 | 84 | 88.4% | 0.25% | 9 | 9 | 98.3% | 0.12% | 4 | 4 | 90.2% | 0.39% |
| | | NOD | 84 | 84 | 82.9% | 0.14% | 9 | 9 | 97.6% | 0.16% | 4 | 4 | 93.1% | 0.18% |
| | | PWD | 92 | 92 | 82.9% | 0.32% | 10 | 10 | 97.2% | 0.19% | 4 | 4 | 88.8% | 0.45% |
| * number of germline genes from reference | | | | | | | | | | | | | | |
| ^ number of germline genes from reference with 0 mismatches | | | | | | | | | | | | | | |
| # percent of reads assigned to unambigious germline reference with 0 mismatches | | | | | | | | | | | | | | |

677

678

679

680

681

682

683

684

685

686

687

688

689

690    **ACKNOWLEDGMENTS**

691

692    **CONFLICTS OF INTEREST**

693    The authors do not have any conflicts of interest to declare.

694

695    **REFERENCES**
696
697

698    1.    Morse HC. Origins of inbred mice. Academic Press: New York, 1978.

699    2.    Alt FW, Baltimore D. Joining of immunoglobulin heavy chain gene segments:

700          implications from a chromosome with evidence of three D-JH fusions. *Proc Natl Acad*

701          *Sci USA* 1982; **79**: 4118-4122.

702    3.    Leder P, Honjo T, Packman S, Swan D, Nau M, Norman B. The organization and

703          diversity of immunoglobulin genes. *Proc Natl Acad Sci USA* 1974; **71**: 5109-5115.

704    4.    Potter M. Antigen-binding myeloma proteins of mice. *Adv Immunol* 1977; **25**: 141-211.

705    5.    Collins AM, Wang Y, Roskin KM, Marquis CP, Jackson KJ. The mouse antibody heavy

706          chain repertoire is germline-focused and highly variable between inbred strains. *Philos*

707          *Trans R Soc Lond B Biol Sci* 2015; **370**.

708    6.    Johnston CM, Wood AL, Bolland DJ, Corcoran AE. Complete sequence assembly and

709          characterization of the C57BL/6 mouse Ig heavy chain V region. *J Immunol* 2006; **176**:

710          4221-4234.

711    7.    Riblet R. Immunoglobulin heavy chain genes in the mouse. In: Honjo T, Alt FW,

712          Neuberger M (eds). Molecular biology of B cells. London: Elsevier Academic Press;

713          2004, pp 19-26.

714    8.    Retter I, Chevillard C, Scharfe M*, et al.* Sequence and characterization of the Ig heavy

715          chain constant and partial variable region of the mouse strain 129S1. *J Immunol* 2007;

716          **179**: 2419-2427.

717    9.    Retter I, Althaus HH, Munch R, Muller W. VBASE2, an integrative V gene database.

718          *Nucleic Acids Res* 2005; **33**: D671-674.

719    10.   Giudicelli V, Duroux P, Ginestoux C*, et al.* IMGT/LIGM-DB, the IMGT comprehensive

720          database of immunoglobulin and T cell receptor nucleotide sequences. *Nucleic Acids*

721          *Res* 2006; **34**: D781-784.

722    11.   Matthyssens G, Rabbitts TH. Structure and multiplicity of genes for the human

723          immunoglobulin heavy chain variable region. *Proc Natl Acad Sci U S A* 1980; **77**: 6561-

724          6565.

725    12.   Matsuda F, Ishii K, Bourvagnet P*, et al.* The complete nucleotide sequence of the human

726          immunoglobulin heavy chain variable region locus. *J Exp Med* 1998; **188**: 2151-2162.

727    13.   Boyd SD, Gaeta BA, Jackson KJ*, et al.* Individual variation in the germline Ig gene

728          repertoire inferred from variable region gene rearrangements. *J Immunol* **184**: 6986-

729          6992.

730    14.   Glanville J, Zhai W, Berka J*, et al.* Precise determination of the diversity of a

731          combinatorial antibody library gives insight into the human immunoglobulin repertoire.

732          *Proc Natl Acad Sci USA* 2009; **106**: 20216-20221.

733  15.  Gidoni M, Snir O, Peres A, *et al.* Mosaic deletion patterns of the human antibody heavy

734      chain gene locus shown by Bayesian haplotyping. *Nature communications* 2019; **10**:

735      628.

736  16.  Scheepers C, Shrestha RK, Lambson BE, *et al.* Ability to develop broadly neutralizing

737      HIV-1 antibodies is not restricted by the germline Ig gene repertoire. *J Immunol* 2015;

738      **194**: 4371-4378.

739  17.  Wang Y, Jackson KJ, Gaeta B, *et al.* Genomic screening by 454 pyrosequencing

740      identifies a new human IGHV gene and sixteen other new IGHV allelic variants.

741      *Immunogenetics* 2011; **63**: 259-265.

742  18.  Watson CT, Steinberg KM, Huddleston J, *et al.* Complete haplotype sequence of the

743      human immunoglobulin heavy-chain variable, diversity, and joining genes and

744      characterization of allelic and copy-number variation. *Am J Hum Genet* 2013; **92**: 530-

745      546.

746  19.  Avnir Y, Watson CT, Glanville J, *et al.* IGHV1-69 polymorphism modulates anti-influenza

747      antibody repertoires, correlates with IGHV utilization shifts and varies by ethnicity. *Sci*

748      *Rep* 2016; **6**: 20842.

749  20.  Watson CT, Glanville J, Marasco WA. The Individual and population genetics of antibody

750      immunity. *Trends Immunol* 2017; **38**: 459-470.

751  21.  Ohlin M, Scheepers C, Corcoran M, *et al.* Inferred allelic variants of immunoglobulin

752      receptor genes: a system for their evaluation, documentation, and naming. *Front*

753      *Immunol* 2019; **10**: 435.

754  22.  Corcoran MM, Phad GE, Vazquez Bernat N, *et al.* Production of individualized V gene

755      databases reveals high levels of immunoglobulin genetic diversity. *Nature*

756      *communications* 2016; **7**: 13642.

757   23.   Gadala-Maria D, Gidoni M, Marquez S, *et al.* Identification of subject-specific

758         immunoglobulin alleles from expressed repertoire sequencing data. *Front Immunol* 2019;

759         **10**: 129.

760   24.   Ralph DK, Matsen IV FA. Per sample immunoglobulin germline inference from B cell

761         receptor deep sequencing data.        *arXiv:171105843v2 [q-bioPE]* 2018.

762   25.   Zhang W, Wang IM, Wang C, *et al.* IMPre: An accurate and efficient software for

763         prediction of T- and B-cell receptor germline genes and alleles from rearranged

764         repertoire data. *Front Immunol* 2016; **7**: 457.

765   26.   Jouvin-Marche E, Morgado MG, Leguern C, Voegtle D, Bonhomme F, Cazenave PA.

766         The mouse Igh-1a and Igh-1b H chain constant regions are derived from two distinct

767         isotypic genes. *Immunogenetics* 1989; **29**: 92-97.

768   27.   Collins AM, Jackson KJL. On being the right size: antibody repertoire formation in the

769         mouse and human. *Immunogenetics* 2018; **70**: 143-158.

770   28.   Yang H, Wang JR, Didion JP, *et al.* Subspecific origin and haplotype diversity in the

771         laboratory mouse. *Nat Genet* 2011; **43**: 648-655.

772   29.   Thornqvist L, Ohlin M. Critical steps for computational inference of the 3'-end of novel

773         alleles of immunoglobulin heavy chain variable genes - illustrated by an allele of IGHV3-

774         7. *Mol Immunol* 2018; **103**: 1-6.

775   30.   Feeney AJ. Predominance of VH-D-JH junctions occurring at sites of short sequence

776         homology results in limited junctional diversity in neonatal antibodies. *J Immunol* 1992;

777         **149**: 222-229.

778   31.   Makino S, Kunimoto K, Muraoka Y, Mizushima Y, Katagiri K, Tochino Y. Breeding of a

779         non-obese, diabetic strain of mice. *Jikken Dobutsu* 1980; **29**: 1-13.

780   32.   Mullen Y. Development of the Nonobese diabetic mouse and contribution of animal

781         models for understanding type 1 diabetes. *Pancreas* 2017; **46**: 455-466.

782    33.    Beck JA, Lloyd S, Hafezparast M*, et al.* Genealogies of mouse inbred strains. *Nat Genet*

783           2000; **24**: 23-25.

784    34.    Staats J. The laboratory mouse. In: Green EL (ed). Biology of the laboratory mouse.

785           New York: McGraw-Hill; 1966, pp 1-9.

786    35.    Prochazka M, Serreze DV, Frankel WN, Leiter EH. NOR/Lt mice: MHC-matched

787           diabetes-resistant control strain for NOD mice. *Diabetes* 1992; **41**: 98-106.

788    36.    Suzuki H, Nunome M, Kinoshita G*, et al.* Evolutionary and dispersal history of Eurasian

789           house mice Mus musculus clarified by more extensive geographic sampling of

790           mitochondrial DNA. *Heredity (Edinb)* 2013; **111**: 375-390.

791    37.    Suzuki H, Yakimenko LV, Usuda D, Frisman LV. Tracing the eastward dispersal of the

792           house mouse, Mus musculus. *Genes and environment* 2015; **37**: 20.

793    38.    Dhande IS, Cranford SM, Zhu Y*, et al.* Susceptibility to hypertensive renal disease in the

794           Spontaneously Hypertensive rat is influenced by 2 loci affecting blood pressure and

795           immunoglobulin repertoire. *Hypertension* 2018; **71**: 700-708.

796    39.    Caron J, Loredo-Osti JC, Laroche L, Skamene E, Morgan K, Malo D. Identification of

797           genetic loci controlling bacterial clearance in experimental Salmonella enteritidis

798           infection: an unexpected role of Nramp1 (Slc11a1) in the persistence of infection in mice.

799           *Genes Immun* 2002; **3**: 196-204.

800    40.    Swihart K, Fruth U, Messmer N*, et al.* Mice from a genetically resistant background

801           lacking the interferon gamma receptor are susceptible to infection with Leishmania major

802           but mount a polarized T helper cell 1-type CD4+ T cell response. *J Exp Med* 1995; **181**:

803           961-971.

804    41.    Hannestad K, Scott H. The MHC haplotype H2b converts two pure nonlupus mouse

805           strains to producers of antinuclear antibodies. *J Immunol* 2009; **183**: 3542-3550.

806   42.   Collins AM, Watson CT. Immunoglobulin light chain gene rearrangements, receptor

807        editing and the development of a self-tolerant antibody repertoire. *Front Immunol* 2018;

808        **9**: 2249.

809   43.   Vander Heiden JA, Yaari G, Uduman M*, et al.* pRESTO: a toolkit for processing high-

810        throughput sequencing raw reads of lymphocyte receptor repertoires. *Bioinformatics*

811        2014; **30**: 1930-1932.

812   44.   Alamyar E, Giudicelli V, Duroux P, Lefranc MP. Antibody V and C domain sequence,

813        structure, and interaction analysis with special reference to IMGT. *Methods Mol Biol*

814        2014; **1131**: 337-381.

815   45.   Gupta NT, Vander Heiden JA, Uduman M, Gadala-Maria D, Yaari G, Kleinstein SH.

816        Change-O: a toolkit for analyzing large-scale B cell immunoglobulin repertoire

817        sequencing data. *Bioinformatics* 2015; **31**: 3356-3358.

818   46.   Huang Y, Niu B, Gao Y, Fu L, Li W. CD-HIT Suite: a web server for clustering and

819        comparing biological sequences. *Bioinformatics* 2010; **26**: 680-682.

820   47.   Meyerhans A, Vartanian JP, Wain-Hobson S. DNA recombination during PCR. *Nucleic*

821        *Acids Res* 1990; **18**: 1687-1691.

822   48.   Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.

823        *J Mol Biol* 1990; **215**: 403-410.

824   49.   Kent WJ. BLAT--the BLAST-like alignment tool. *Genome Res* 2002; **12**: 656-664.

825   50.   Ye J, Ma N, Madden TL, Ostell JM. IgBLAST: an immunoglobulin variable domain

826        sequence analysis tool. *Nucleic Acids Res* 2013.

827

828   **REFERENCES**

829   **Additional References for Supplementary figures**

830        1.    Wang JR, de Villena FP, McMillan L. Comparative analysis and

831     visualization of multiple collinear genomes. *BMC Bioinformatics* 2012; **13:** S13.

832     2.      Yang H, Wang JR, Didion JP *et al.* Subspecific origin and haplotype

833     diversity in the laboratory mouse. *Nat Genet* 2011; **43:** 648-655.

834

835

836     **FIGURE LEGENDS**

837     **Figure 1.** Comparisons of inferred germline IGHV sequences from each strain

838     (CAST/EiJ, n=1; LEWES/EiJ, n=1; MSM/MsJ, n=1; PWD/PhJ, n=1; NOD/ShiLtJ, n=1) to

839     those represented in the mouse IMGT database (imgt.org). (**a**) Donut plots depicting the

840     proportion of inferred germline sequences from each strain that align to a known IMGT

841     allele with 100% match identity. (**b**) Boxplots depicting the sequence similarities of

842     inferred germline sequences from each strain when compared to the closest known

843     IMGT allele. (**c**) The count of identified germline sequences from each strain

844     representing known mouse IGHV gene subfamilies. (**d**) The count of non-IMGT inferred

845     germline sequences in each strain, partitioned by IGHV gene subfamily.

846

847     **Figure 2.** IGHD and IGHJ inferred germline sequences among mouse strains

848     (CAST/EiJ, n=1; LEWES/EiJ, n=1; MSM/MsJ, n=1; PWD/PhJ, n=1; NOD/ShiLtJ, n=1).

849     (**a**) Tile plot depicting the presence of IGHD genes across the different mouse strain.

850     Filled black cells indicate the presence of the gene. For novel variants, the sequences

851     are listed along with the *01 reference allele, with nucleotide differences for the variant

852     indicated in uppercase bold. (**b**) Tile plot depicting the presence of IGHJ genes across

853     the different mouse strains. Black cells indicate that a gene was confirmed as present in

854     a strain, white indicates that a gene wasn't confirmed for a strain and gray indicates

37

855      genes whose presence remains uncertain. For the novel variant, the sequence is listed

856      along with the *01 reference allele, with nucleotide differences for the variant indicated

857      in uppercase bold.

858

859      **Figure 3.** Relationships of IGHV inferred germline sequences among mouse strains

860      (CAST/EiJ, n=1; LEWES/EiJ, n=1; MSM/MsJ, n=1; PWD/PhJ, n=1; NOD/ShiLtJ, n=1).

861      (**a**) Upset plot depicting the size of the germline set from each of the analyzed strains

862      (left), as well as the numbers of sequences either unique to a given strain or shared

863      among strains (identical sequences). (**b**) Heatmap depicting the mean percent

864      sequence match identities among inferred IGHV germline sets for each pair-wise strain

865      comparison (see also Supplementary figure 4).

866

867      **Supplementary figure 1.** Single nucleotide variant (SNV) data from the IGHV gene

868      region (chr12:114700000-117270000) suggest common subspecies origins for IGHV

869      genomic haplotypes of inbred and wild-derived laboratory mouse strains. This figure

870      depicts the predicted subspecies origins (*Mus musculus domesticus*; *M. m. musculus*;

871      *M. m. castaneous*) of IGHV haplotypes in the six strains analyzed in the present study.

872      Here, we consider the relationships between inferred germline IGHV gene sets of these

873      strains in the context of these predicted subspecific origins. Data presented in this figure

874      were obtained from the Mouse Phylogeny Viewer [1] (https://msub.csbio.unc.edu/, based

875      on previously published whole-genome SNV data [2]. The bars above each haplotype

876      depict the locations of "diagnostic" SNVs used to make subspecies determinations.

877

878 **Supplementary figure 2.** Counts of identified clones representing inferred germline

879 sequences from each strain sequenced in this study. Additional information, including

880 the nucleotide sequences of each inferred germline presented in these plots, can be

881 found in Supplemental Table 2.

882

883 **Supplementary figure 3.** Boxplots depicting the numbers of clones representing each

884 inferred IGHV germline sequence in CAST/EiJ, LEWES/EiJ, MSM/MsJ, and PWD/PhJ,

885 respectively. Data for each strain are partitioned by whether the inferred sequence was

886 supported by IgDiscover. With only a few exceptions, the inferred IGHV sequences in

887 each strain lacking support by IgDiscover ("No") were represented by fewer clones in

888 our analysis on average than those that were supported by IgDiscover ("Yes").

889

890 **Supplementary figure 4.** Violin plots depicting the best percent match identities among

891 inferred IGHV germline sets for each pair-wise strain comparison. Individual points

892 within each violin plot represent the best percent identity value of an IGHV gene

893 segment as compared to all genes in the other strain to which it was compared.

894

895 **Supplementary table 1.** Per-strain AIRR-seq library summary statistics for CAST/EiJ,

896 LEWES/EiJ, MSM/MsJ, PWD/PhJ, and NOD/ShiLtJ.

897 **Supplementary table 2.** Complete database of inferred germline IGHV sequences from

898 CAST/EiJ, LEWES/EiJ, MSM/MsJ, PWD/PhJ, and NOD/ShiLtJ.

899 **Supplementary table 3.** Complete dataset of inferred germline IGHD sequences from

900 CAST/EiJ, LEWES/EiJ, MSM/MsJ, PWD/PhJ, and NOD/ShiLtJ.

901     **Supplementary table 4.** Complete dataset of inferred germline IGHJ sequences from

902     CAST/EiJ, LEWES/EiJ, MSM/MsJ, PWD/PhJ, and NOD/ShiLtJ.
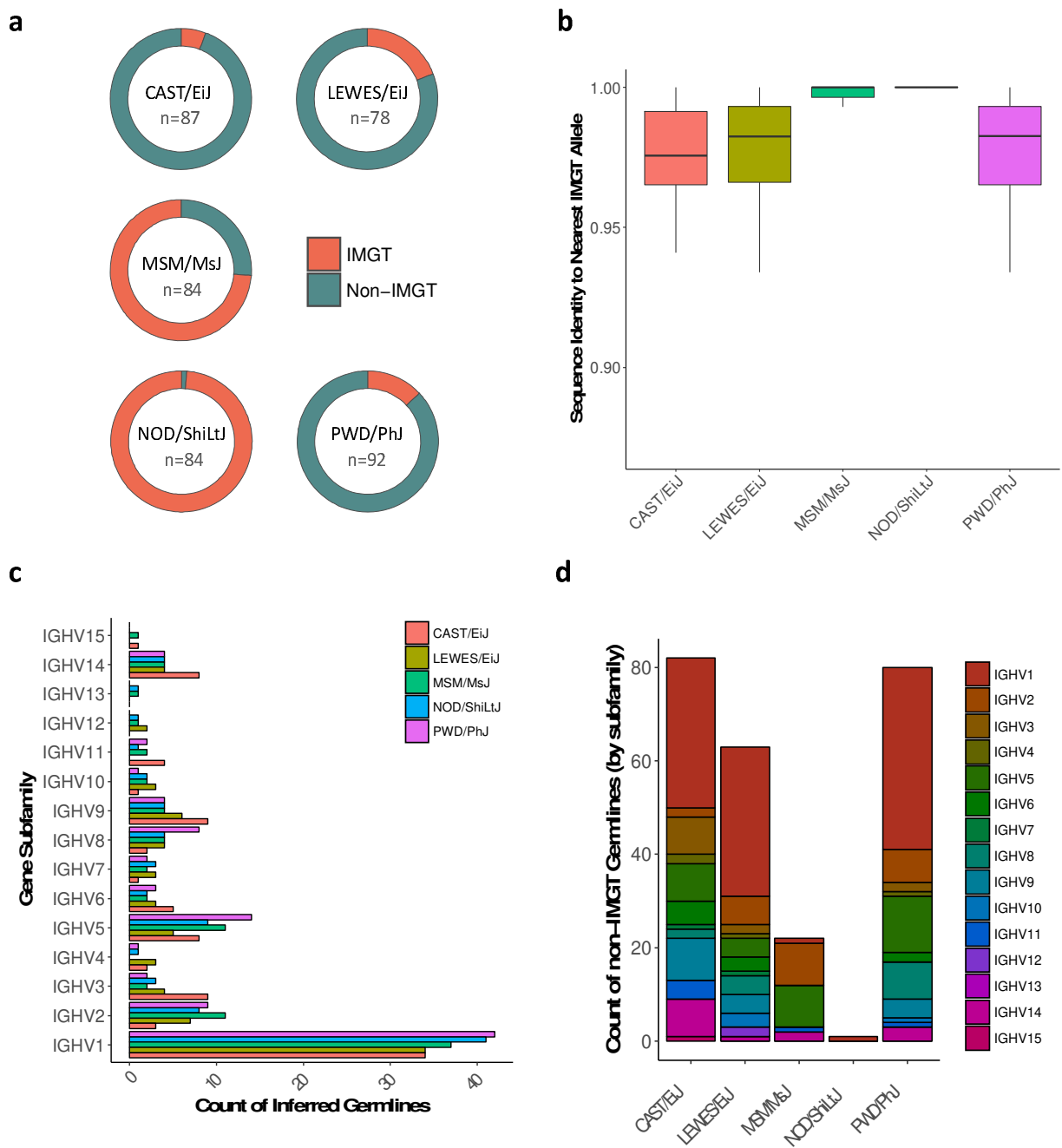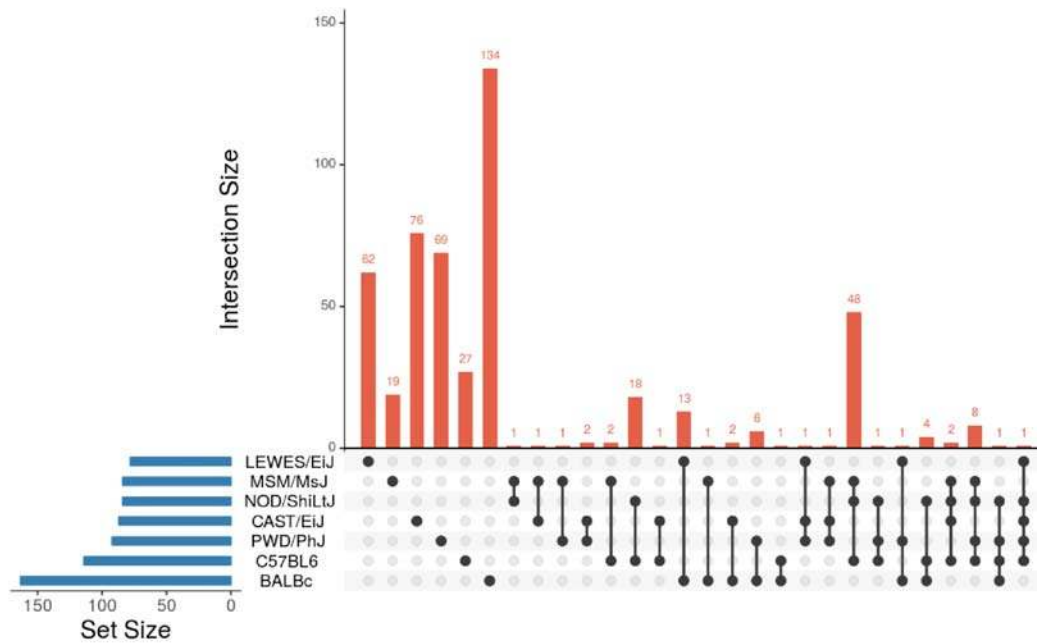
903

904

Figure 1

Figure 2

a



| | BALBc | C57BL/6 | NOD | LEWES | PWD | MSM | CAST | GENE SEQUENCE |
|---|---|---|---|---|---|---|---|---|
| IGHD1−1*01 | | | | | | | | tttattactacggtagtagctac |
| IGHD1−1_varA | | | | | | | | tttattactac**Agt**G**g**tagctac |
| IGHD1−1_varB | | | | | | | | tttattactacg**At**G**g**tagctac |
| IGHD1−2*01 | | | | | | | | tttattactacg**At**G**g**tagctac |
| IGHD1−3*01 | | | | | | | | |
| IGHD2−1*01 I IGHD2−8*01 | | | | | | | | |
| IGHD2−10*01 I IGHD2−11*02 | | | | | | | | |
| IGHD2−10*02 I IGHD2−11*01 | | | | | | | | |
| IGHD2−12*01 | | | | | | | | tctatgatggttactac |
| IGHD2−12_var | | | | | | | | tctatg**G**tggttactac |
| IGHD2−13*01 | | | | | | | | |
| IGHD2−14*01 | | | | | | | | |
| IGHD2−2*01 I IGHD2−7*01 | | | | | | | | |
| IGHD2−3*01 | | | | | | | | tttattactacggtagtagctac |
| IGHD2−3_var | | | | | | | | tttattactac**Agt**G**g**tagctac |
| IGHD2−4*01 I IGHD2−9*02 | | | | | | | | |
| IGHD2−5*01 I IGHD2−6*01 | | | | | | | | |
| IGHD2−9*01 | | | | | | | | |
| IGHD3−1*01 | | | | | | | | |
| IGHD3−2*01 | | | | | | | | |
| IGHD3−2*02 | | | | | | | | |
| IGHD3−3*01 | | | | | | | | |
| IGHD4−1*01 | | | | | | | | |
| IGHD4−1*02 | | | | | | | | |

b



| | BALBc | C57BL/6 | NOD | LEWES | PWD | MSM | CAST | GENE SEQUENCE |
|---|---|---|---|---|---|---|---|---|
| IGHJ1*01 | | | | | | | | |
| IGHJ1*03 | | | | | | | | |
| IGHJ2*01 | | | | | | | | actactttgactactggggccaaggcaccactctcacagtctcctcag |
| IGHJ2_var | | | | | | | | actactttgactactggggccaaggcaccact**At**cacagtctcctcag |
| IGHJ3*01 | | | | | | | | |
| IGHJ4*01 | | | | | | | | |

Figure 3

**a**



**b**