1  Characterisation of the immune repertoire of a humanised transgenic mouse

2  through immunophenotyping and high-throughput sequencing

3

4  Richardson E[1,2*], Binter Š[1*], Kosmac M[1], Ghraichy M[3,4], von Niederhausern V[3,4], Kovaltsuk A[2], Galson J[5],

5  Trück J[3,4], Kelly DF[6], Deane CM[2], Kellam P[1, 7], **Watson SJ**[1]

6

7  [1] Kymab, a Sanofi Company, Babraham Research Campus, Cambridge, UK

8  [2] Department of Statistics, University of Oxford, Oxford, UK

9  [3] Division of Immunology, University Children's Hospital, University of Zurich, Zurich, Switzerland

10  [4] Children's Research Center, University of Zurich, Zurich, Switzerland

11  [5] Alchemab Therapeutics Ltd, Kings Cross, London, UK

12  [6] Department of Paediatrics, University of Oxford, Oxford, UK

13  [7] Department of Infectious Diseases, Faculty of Medicine, Imperial College London, UK

14

15  [*] Contributed equally to the study

16  **Corresponding author: Simon Watson (simon.watson@sanofi.com)**

17

18

# Abstract

19

20   Immunoglobulin loci-transgenic animals are widely used in antibody discovery and increasingly in

21   vaccine response modelling. In this study, we phenotypically characterised B-cell populations from the

22   Intelliselect® Transgenic mouse (Kymouse) demonstrating full B-cell development competence.

23   Comparison of the naïve B-cell receptor (BCR) repertoires of Kymice BCRs naïve human and murine

24   BCR repertoires revealed key differences in germline gene usage and junctional diversification. These

25   differences result in Kymice having CDRH3 length and diversity intermediate between mice and

26   humans. To compare the structural space explored by CDRH3s in each species repertoire, we used

27   computational structure prediction to show that Kymouse naïve BCR repertoires are more human-like

28   than mouse-like in their predicted distribution of CDRH3 shape. Our combined sequence and

29   structural analysis indicates that the naïve Kymouse BCR repertoire is diverse with key similarities to

30   human repertoires, while immunophenotyping confirms that selected naïve B-cells are able to go

31   through complete development.

32   **Word count: 147**

## Introduction

Twenty-five years of progress in genetic engineering since the first Ig transgenic mouse (M. Brüggemann et al. 1989) culminated in 2014 in the insertion of the full set of variable human Ig genes in mice (Lee et al. 2014). Humanised immunoglobulin (Ig) loci-transgenic animal models have proven extremely useful in therapeutic antibody discovery and, more recently, in vaccine response modelling; twenty of the 127 therapeutic antibodies licensed in the US or EU as of April 2022 were derived from transgenic mouse platforms (data from Thera-SAbDab (Raybould et al. 2020)). Transgenic platforms have also found a new application in vaccine response modelling (Sok et al. 2016; Pantophlet et al. 2017; Walls et al. 2020). As humanised animal models become the source of a growing number of therapeutics and play an increasingly important role in the evaluation of novel vaccine candidates, it is crucial to understand the degree to which their B-cell repertoires can be considered representative of humans.

Contemporary Ig transgenic animal models vary according to the number of genes and localisation of the inserted human Ig loci (Green 2014; Brüggemann et al. 2015). In Kymab's Intelliselect® Transgenic mouse (Kymouse), a complete set of human variable (V), diversity (D) and junction (J) genes of the IGH locus as well as the V and J genes of the IGλ and IGκ loci were inserted at the sites of the endogenous mouse loci. The mouse constant regions were retained, preserving downstream interactions with endogenous intracellular signalling components and cell membrane Fc receptors, resulting in functional, fully active chimeric antibodies. Kymice exhibit normal B-cell production and maturation and the resulting B-cell receptors (BCRs) are diverse, with human-like CDRH3 lengths and evidence of somatic hypermutation (Lee et al. 2014). However, the baseline phenotypic diversity in B-cells and B-cell receptors (BCRs) in the Kymouse has not been fully described.

B cells are an integral part of the humoral immune response due to their ability to produce antibodies against diverse antigens, providing protection against infection. They originate from hematopoietic stem cells in the bone marrow, where they undergo several phases of antigen-independent

3

58  development leading to the generation of immature B cells. B cells are routinely classified based on

59  their maturation status, antibody isotype, and effector function. Ig gene rearrangement during these

60  early stages of B cell development results in the expression of a mature B cell receptor that is capable

61  of binding to antigen. This is followed by positive and negative selection processes, which are designed

62  to eliminate non-functional and self-reactive immature B cells. Surviving B cells complete antigen-

63  independent maturation in the spleen, producing immunocompetent naïve mature B cells that

64  subsequently develop into either follicular or marginal zone B cells. In response to vaccination or

65  invading microbes, antigen-specific B cells within secondary lymphoid organs differentiate into

66  antibody-producing cells, early memory cells, or rapidly proliferate and form structures known as

67  germinal centres (Allen, Okada, and Cyster 2007). Germinal centres (GCs) are inducible lymphoid

68  microenvironments that support the generation of affinity-matured, isotype-switched memory B cells

69  and antibody secreting plasma cells. Long-lived plasma cells (PCs) secrete high-affinity antibodies, and

70  memory cells can readily elicit an efficient antibody immune response upon re-exposure to the

71  immune stimuli (Corcoran and Tarlinton 2016; Weisel and Shlomchik 2017). Iterative cycles of B cell

72  hypermutation and selection within the GC leads to an accumulation of affinity-enhancing mutations

73  and ultimately to the progressive increase of serum antibody affinity, a process known as antibody

74  affinity maturation (Jacob et al. 1991). Antibody-secreting plasma cells play critical roles in protective

75  immunity on the one hand and antibody mediated autoimmune disease on the other. During immune

76  responses a small fraction of newly generated plasma cells enter either the bone marrow (BM) or the

77  lamina propria of the small intestine where they populate specialised survival niches and become long-

78  lived plasma cells (Lemke et al. 2016) thus maintaining antibody titres for extended periods.

79  The variable domain of a BCR is composed of a heavy chain and a light chain. Each of the chains in the

80  antibody has three hypervariable regions known as the complementarity-determining regions (CDRs),

81  which make most contacts with the antigen. The heavy chain locus consists of variable (V), diversity

82  (D) and joining (J) gene segments, which recombine to form the heavy chain: the first two CDRs of the

83  heavy chain, CDRH1 and CDRH2, are encoded by the V gene alone, while the third and most variable

4

84    CDR, CDRH3, spans the V, D and J gene junctions. The insertion of random and palindromic nucleotides

85    at the V-D and D-J junctions further contribute to the diversity of the CDRH3, ensuring binding diversity

86    to different antigens and epitopes (Xu and Davis 2000). Each of the light chain loci, kappa and lambda,

87    consist of V and J gene segments but no D gene segments, and both the germline as well as the

88    recombined chain are less diverse than their heavy chain counterparts (Collins and Watson 2018).

89    Due to the greater diversity of the heavy chain, most next-generation sequencing of BCR repertoires

90    (BCR-seq) has focused on the heavy chain; lower throughput methods exist for identifying the light

91    chain pairing (Curtis and Lee 2020). The resulting BCR sequences can be aligned to reference germline

92    gene databases to infer most likely germline gene origins and insertions or deletions at the V(D) or

93    (D)J junctions (Ye et al. 2013). Alignment of BCRs to common germline genes also allows inference

94    about clonal structure, as sequences sharing common germline gene assignments as well as homology

95    in the CDRH3 loop may be inferred to have arisen from a common progenitor B-cell (Greiff et al. 2015;

96    Yaari and Kleinstein 2015). The amino acid sequences of these heavy chains can also be functionally

97    examined through annotation with structural tools (Kovaltsuk et al. 2017; Marks and Deane 2020).

98    Changes in the pattern of CDRH3 shapes in BCR repertoires have been observed along the B-cell

99    differentiation axis in both humans and mice (Kovaltsuk et al. 2020) but the extent to which the CDRH3

100   protein shape differs between humans and mice has not been explored.

101   Here we have characterised the frequency of germinal centre B cells, memory B cells and long-lived

102   plasma cells from spleens, lymph nodes, and bone marrows of antigen-inexperienced Kymice (Lee et

103   al. 2014). The frequencies of these B cell subsets as well as the breadth and nature of their BCR

104   repertoires constitutes the first step in our understanding of how the immune system of this model

105   organism responds to different antigens, vaccines, and pathogens that are both of scientific as well as

106   therapeutic interest. Examining the nature of the naïve BCR repertoire in Kymice through both single-

107   cell and bulk sequencing and structural analysis shows the Kymouse naïve BCR repertoires are more

108   human-like in their distribution of CDRH3 shapes.

109

# Materials & Methods

**Kymouse data**

112  Spleens, lymph nodes and bone marrows were collected from *n*=20 antigen inexperienced

113  Intelliselect® Transgenic mice (Kymice). These Kymice contain chimeric immunoglobulin loci, with

114  humanised variable domains ($V_H$, $V_K$, and $V_L$) and a humanised lambda constant domain ($C_L$), but

115  murine heavy ($C_H$) and kappa ($C_K$) constant domains (Lee et al. 2014). The Kymice were selected to

116  reduce any possible confounding effects by ensuring that: (i) there was an equal representation of

117  sexes, (ii) mice were a range of ages on culling (6-12 weeks old), and (iii) mice were selected from

118  different litters and culled over a period of 6 months.

119

**Lymphoid Cell Isolation and Cryopreservation**

121  Bone marrow isolated from the femurs and tibias of each Kymouse were processed to single-cell

122  suspensions by flushing the tissues with ice-cold FBS buffer and passing through 40 µm cell strainers.

123  Spleens and inguinal lymph nodes were processed to single cell suspensions by homogenising through

124  40 µm cell strainers with ice-cold FBS buffer and pooled prior to staining and cell sorting. All single-

125  cell suspensions were pelleted at 400 g for 10 minutes at 4°C prior to cryopreservation in 10%

126  DMSO/FBS and storage in liquid nitrogen.

127

**Next Generation Sequencing (NGS) Analysis of Paired $V_H$ and $V_L$ Sequences from Single-Cell**

**Sorted B Cells Derived from Kymice**

130  For each Kymouse, the spleen and inguinal lymph nodes were processed to single-cell suspensions as

131  described above before fluorescently-activated cell sorting (FACS) to recover $CD19^+$ $B220^+$ B cells into

132  individual wells of a 96 well plate. RT-PCR was performed to amplify the $V_H$ and $V_L$ domains, and

133  standard Illumina libraries were generated before sequencing on an Illumina MiSeq sequencer. The

134  Change-O pipeline (Gupta et al. 2015) was used to process the sequence data; naïve BCR sequences

135  were characterized as immature B-cells with sequences containing zero nucleotide mutations. In total,

136  $n$=3,885 paired $V_H$ and $V_L$ sequences were processed from the $n$=20 Kymice.

137

138  **NGS Sequence Analysis of $V_H$ Sequences from Bulk Sorted B Cells Derived from Kymice**

139  Bone marrows from the femur and tibia of each Kymouse were processed into single-cell suspension

140  as described above. From these $n$=20 bone marrow samples, seven were FAC sorted to recover $CD19^+$

141  $B220^+$ B cells into a single tube. The cells were lysed and RT-PCR was performed to amplify the $V_H$

142  domain, followed by standard Illumina library generation, before sequencing on an Illumina MiSeq.

143  The Change-O pipeline (Gupta et al. 2015) was used to process the sequences generated by the MiSeq

144  sequencers. IgM sequences with zero mutations were selected for further analysis resulting in a total

145  of 412,493 $V_H$ sequences across the n = 7 Kymice (average 58,928 range 31,905 to 100,240).

146

147  **NGS Sequence Analysis of $V_H$ Sequences Derived from Human Samples**

148  Buffy coat samples were obtained from ten healthy individuals as described previously (Ghraichy et

149  al, 2021). In the previous study, B-cells were FAC-sorted into naïve, marginal zone, plasma and

150  switched memory cell populations. In the present study, we analysed IgM sequences from the naïve

151  subset of B-cells with no mutations. There was a total of $n$=338,677 sequences (mean: 33,867 per

152  human, range: 20,653 – 48,293).

153

7

154 **NGS Sequence Analysis of V$_H$ Sequences Derived from C57BL/6 WT Mice**

155 6,763,480 IgM V$_H$ nucleotide sequences from naïve B cells of healthy unvaccinated C57BL/6 wild-type

156 mice were downloaded from the Observed Antibody Space (OAS) (Kovaltsuk et al. 2018; Olsen, Boyles,

157 and Deane 2022). Those sequences with any nucleotide mutations were removed, and the remaining

158 sequences were down-sampled via stratified sampling aiming to preserve the original clonal structure

159 in the complete dataset. 150,000 sequences with redundancy were randomly selected from each of

160 the five C57BL/6 mice. Collapsing to unique (nucleotide) sequences resulted in a total of 268,285

161 sequences (mean: 53,657 sequences per mouse; range: 20,026-87,041).

162

163 **Clonal and diversity analysis of NGS sequence data**

164 Clonotypes are defined as sequences with common IGHV and IGHJ genes and 90% or more amino acid

165 identity across length-matched CDRH3s. Antibody sequences were assigned to clonotypes using the

166 DefineClones module of Change-O (Gupta et al. 2015) under the amino acid model.

167  Shannon diversity (H) was calculated using the *stats.entropy* function within Python's SciPy library.

168 The formula is as follows:

169
$$H = -\sum_{i=1}^{s} p_i ln p_i$$

170 Where $p_i$ is the proportion of sequences in the clonotype *i* of *s* clonotypes.

171

172 # Structural annotation

173 Sequences that IgBLAST identified as coding for a productive immunoglobulin were translated into

174 amino acids and aligned to the IMGT antibody numbering scheme (Lefranc et al. 2003) using ANARCI

8

175   (Dunbar and Deane 2016). IMGT CDR sequences were extracted and structurally annotated using the

176   SAAB+ version 1.01 pipeline (Kovaltsuk et al. 2020). SAAB+ uses SCALOP (Wong et al. 2019) to assign

177   each sequence's CDRH1 and CDRH2 loops to a structural canonical class and uses FREAD (Choi and

178   Deane 2010) to identify if the CDRH3 loop has a similar structure to any crystallographically-solved set

179   of 4,544 CDRH3 structures (referred to as templates, downloaded from SAbDab (Dunbar et al. 2014;

180   Schneider, Raybould, and Deane 2022) on 16th February 2022).

181   To reduce dimensionality, templates are clustered with a 0.6Å RMSD cut-off, producing a set of 1,944

182   templates. In this set of templates, 41% of the antibody structures are of murine-origin and 37% are

183   of human-origin (**Supplementary Figure 5**). The SAAB+ pipeline outputs for each sequence the

184   canonical class of the CDRH1 and CDRH2 loops, and the Protein Data Bank (PDB) ID of the structure

185   that contains the best matched CDRH3 structure for homology modelling from the 1,944 templates.

186   As a complementary approach, we also modelled representatives of all non-singleton clonotypes in

187   each of the repertoires using a recent deep learning method, ABlooper (Abanades et al. 2022).

188   ABlooper is competitive with AlphaFold2 on the canonical CDRs and shows better performance on the

189   CDRH3 which was the target of our analysis. It is also more than 1000x faster than AlphaFold2 and was

190   therefore suitable for our large-scale analysis.

191   For full structural modelling with ABlooper, a light chain must be supplied. As we did not know the

192   cognate light chain for any of the heavy chains from bulk sequencing, all heavy chains were paired

193   with a single light chain which was selected from the paired dataset. This light chain was the most

194   commonly observed light chain in the Kymouse dataset. We selected a common light chain to

195   standardise its effect on the prediction of the heavy chain CDRs, in the absence of knowledge of the

196   true light chain.

197

198   **Humanness scoring of V$_H$ sequences**

9

199    The human, Kymouse and C57BL/6 mouse $V_H$ sequences were assigned a "humanness score" using

200    the random forest regressors from Hu-mAb (Marks et al. 2021). Sequences were first IMGT numbered

201    using ANARCI as above. While the C57BL/6 mouse and Kymouse sequences were mostly full length

202    (IMGT positions 1-128), the human sequences were in most cases missing FWR1 (IMGT positions 1-26

203    in the IMGT CDR definition). As the human sequences were non-mutated, we considered it reasonable

204    to simply fill in FWR1 according to the sequence found in the assigned germline. For the human and

205    Kymouse sequences, the RF model trained on the IGHV gene assigned by IgBLAST were used for

206    scoring, while for the murine sequences, all seven (IGHV1-7) RF models were used to score the

207    sequence and the highest score was selected. We used the IGHV-specific classification thresholds

208    reported in the Hu-mAb paper to annotate if a sequence was considered human or not (Marks et al.

209    2021).

210

211    **Immunophenotyping of Intelliselect® Transgenic mice**

212    Spleens, lymph nodes, and bone marrow from a further *n*=12 antigen-inexperienced Kymice were

213    processed to single-cell suspensions and cryopreserved as described above. Fluorescently Activated

214    Cell Sorting (FACS) of the bone marrow samples was performed using fluorescently conjugated

215    antibodies against B220, IgM, IgD, IgG1, IgG2ab, IgG3, CD8, CD4, Ly-6G, CD11c and CD138 (BD

216    Biosciences), CD19, F4/80, Sca-1 (BioLegend) and TACI (eBioscience). For the pooled spleen and lymph

217    node samples the FACS panel consisted of B220, IgM, IgD, IgG1, IgG2ab, IgG3, CD8, CD4, Ly-6G&6C,

218    CD11c and CD95 (BD Biosciences), CD19, F4/80, CD73, CD80, PD-L2 and GL7 (BioLegend). DRAQ7

219    (BioStatus) was used in all samples to distinguish live and dead cells. For flow cytometry, cells were

220    thawed from frozen, resuspended in warm 10% FBS in RPMI buffer, filtered through 40 μm cell

221    strainers and centrifuged at 400 × g for 10 minutes at 4 °C. Cells were resuspended in buffer and

222    TrueStain FcX (BioLegend) was added for 10 minutes on ice. Single cell suspensions of bone marrow

223    cells and pooled spleen and lymph node cells were stained with their respective staining cocktails for

224  30 minutes. All cells were spin washed and resuspended in buffer, filtered through a 30 μm cell

225  strainer, followed by cell sorting on a 5-laser BD FACS Aria Fusion flow cytometer (Beckton Dickinson).

226

227  **Analysis of Flow Cytometry Data from Intelliselect® Transgenic mice**

228  The frequencies of the following cell types within total viable (DRAQ7) cells were determined using

229  classical FACS gating: bone marrow plasma cells (CD138$^+$TACI$^+$/Sca-1$^+$), spleen/lymph node memory B

230  cells (B220$^+$CD19$^+$IgD$^-$CD73$^+$/CD80$^+$/PDL2$^+$) and spleen/lymph node germinal centre cells

231  (B220$^+$CD19$^+$CD95$^+$GL7$^+$; data not shown). Cell populations were also analysed in an unbiased manner

232  using unsupervised clustering algorithms. In brief, the raw .fcs files were imported into R (RStudio

233  version 1.2.5033 with R version 4.0.0) using CytoExploreR (version 1.0.8) and the data were

234  transformed to normalise marker intensities (logicle transform). For visualization, additional quantile

235  scaling from 0-1 was performed, fixing values less than the 1$^{st}$ percentile to 0.01 and values above the

236  99$^{th}$ percentile to 0.99 to minimise the contribution of outliers to the scaling. Cell clustering was

237  performed using the Leiden clustering algorithm (R package Monocle 2.16) and clusters were

238  visualized in two-dimensional space using UMAP (R package uwot 0.1.10). Poorly resolved clusters

239  were re-clustered, the subclusters manually merged to the first level clusters and annotated by cell

240  type.

241

242  **Data availability**

243  The sequence data is currently available in the Observed Antibody Space. Immunophenotyping data

244  will be available upon publication.

245

246  **Code availability**

11

247    The Python code used to analyse the data and generate the figures is available at

248    https://github.com/oxpig/HumMus.

249

## Results

**Antigen naïve Kymice exhibit similar B cell sub-population frequencies**

252    We characterised the B cell sub-populations within spleen and lymph node samples of 12 antigen-

253    inexperienced Kymice using an 11-colour flow cytometry panel that incorporated a range of B cell

254    lineage markers to identify both murine memory and germinal centre B cell populations. A canonical

255    gating scheme organises B cells by their maturation status – from transitional B cells through naïve,

256    non-switched and ultimately class-switched memory B cells. To look at the heterogeneity of the B cell

257    subpopulations in more detail we incorporated unbiased Leiden clustering on the multi-parameter

258    FACS data. Sorted cells separated into two large clusters, B-cells and non-B-cells (**Figure 1A**). As

259    expected, within the B cell population immature isotypes (IgD and IgM) were enriched in naïve cells,

260    whereas markers CD95 and GL7 were enriched in germinal centre cells (**Figure 1C**). The murine

261    memory B cell population has been described to comprise five subpopulations defined by the

262    progressive transition from naive-like to more memory-like cells and the surface markers CD80, CD73

263    and PD-L2 have previously been reported to enable their distinction (Tomayko et al. 2010). Using a

264    low dimensional UMAP representation we observed distinct staining patterns of these markers in the

265    memory B cell compartment and were able to distinguish between 12 major B cell populations,

266    including transitional, naïve and activated as well as six distinct memory subsets. These were defined

267    as (1) PD-L2$^{hi}$, (2) CD73$^{hi}$ CD80$^{hi}$ PD-L2$^{low}$, (3) CD80$^{low}$, (4) PD-L2$^{hi}$ CD80$^{hi}$ CD73$^{low}$, (5) CD80$^{hi}$ and (6)

268    CD73$^{hi}$ (**Figure 1A**). Germinal centre B cells formed a small and well separated cluster whose small

269    frequency was not surprising given that these were antigen naïve animals. Based on the 3 memory

270    markers (CD80, CD73 and PD-L2) the relative frequency of the total memory B cells was 6.60% ± 2.51%,

12

271     and the frequency of CD95 and GL7 positive germinal centre cells was 0.18% ± 0.26%. The median

272     expression profile of each subset is shown as a heatmap (**Figure 1C**). The antigen naïve B cell

273     populations in un-immunised and non-infected Kymice are therefore normal and consistent between

274     different Kymice (**Figure 1C right panel**).

275

276     **B cells in the bone marrow are class switched with variable levels of surface BCR expression**

277     To understand the B cell profile beyond spleen and lymph nodes we also profiled the bone marrow

278     cells of mice. We characterized bone marrow samples using a 9-colour flow cytometry panel that

279     incorporated a range of B cell lineage markers. The staining panel was designed to identify plasma

280     cells as well as class switched B cell subsets in antigen inexperienced mice. As expected, we saw that

281     the cells separated again into two large clusters, B-cells and non-B-cells **(Figure 1B).** The expression

282     profiles of the subsets were again plotted as heatmaps showing the median expression profiles of

283     each subset (**Figure 1D**). Within the B cell cluster, we identified several discrete subsets marked by the

284     expression of different BCR isotypes. B cells were clustered into 5 distinct subpopulations, including

285     immature IgD+ and IgM+ cells, mature IgM+ and IgG+ B cells and plasma cells. The markers TACI and

286     Sca-1 were enriched in plasma cells as expected, whereas CD138, a common plasma cells marker did

287     not show a high level of separation between the different cell types. Unsurprisingly, we saw the

288     biggest separation between B cells and non B cells, and a continuum of B cell subtypes from IgD,

289     through IgM, to IgG-expressing B cells as well as a discrete cluster identified as plasma cells. The

290     frequencies of the plasma cells were low (0.90% ± 0.28%) in comparison to other B cell subtypes,

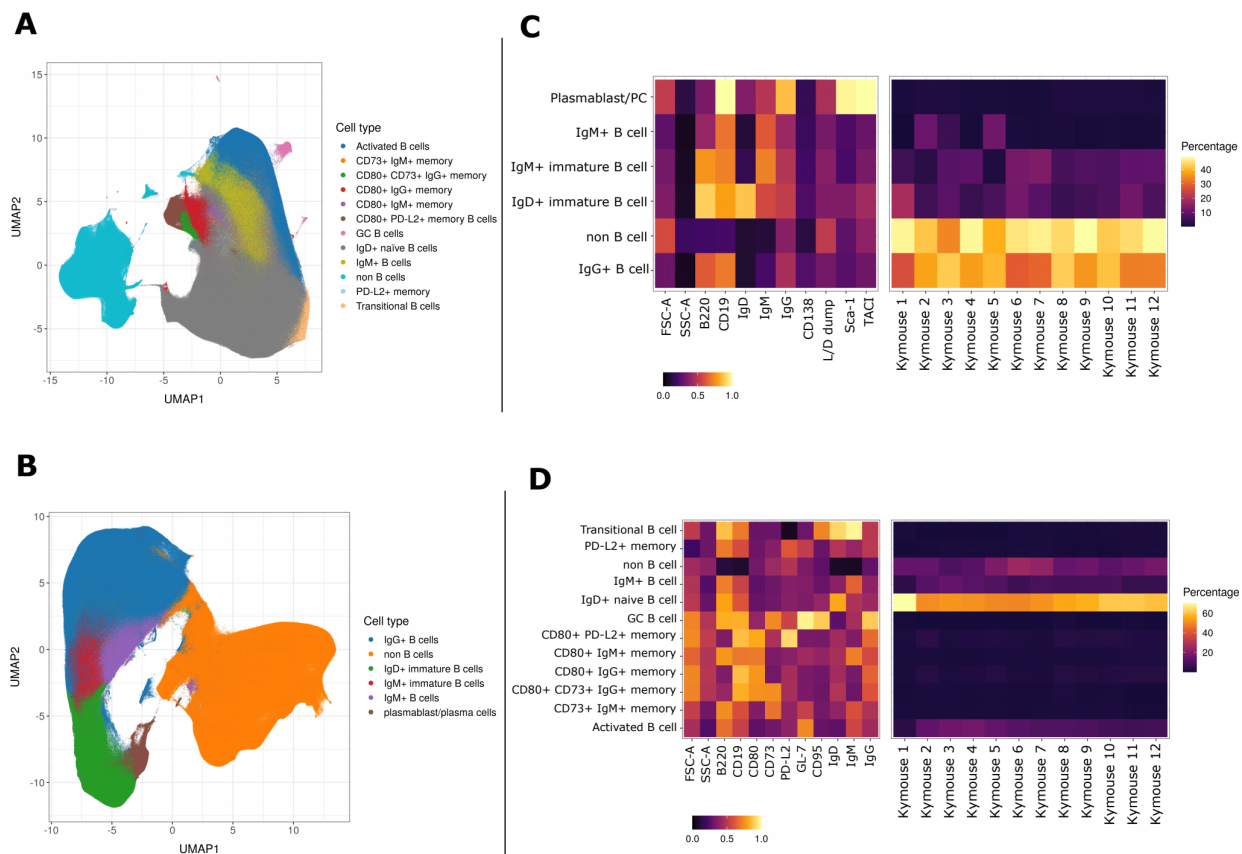291     perhaps not surprising given that these were antigen naïve animals.

*Figure 1:* UMAP projections of sorted cell populations identified using unsupervised clustering from spleen and lymph nodes (A) or from bone marrow samples (B). UMAP projections show a clear separation between B-cells and non-B-cells for both sample types. The projections are coloured by the 12 resolved cell types in the spleen and lymph nodes (A) and the six resolved in the bone marrow samples (B). Normalised and scaled marker expression and frequencies showing mouse-to-mouse variation for each of the resolved cell types in the spleen and lymph nodes (C) or the bone marrow samples (D). The expression profiles are homogeneous across mice.

**The Kymouse naïve antibody sequence repertoire is more human-like than murine-like**

Using high-throughput sequencing we recovered 3,885 full-length paired $V_H$ and $V_L$ sequences and a further 451,655 full-length unpaired $V_H$ sequences from naïve B cells extracted from the spleens and lymph nodes of 20 Kymice. In order to evaluate the humanness of the Kymouse naïve B cell sequence

14

305    repertoire, we compared these sequences to equivalent datasets of 338,677 $V_H$ sequences from

306    human naive B cells, and 268,285 $V_H$ sequences from C57BL/6 mice.

307    One of the most pronounced differences in heavy/light chain pairing between wild type mice and

308    humans that has been described is the usage ratio of the Ig$\kappa$ and Igλ chains in the BCRs of circulating

309    B cells. Humans have an Ig$\kappa$/Igλ ratio of approximately 60:40 in serum and in mature B cells. Mice

310    have an Ig$\kappa$/Igλ ratio of 95:5 in serum and 90:10 on B cells (McGuire and Vitetta 1981). We used the

311    3,885 full-length paired $V_H$ and $V_L$ sequences to calculate the Ig$\kappa$/Igλ ratio in Kymice and found a ratio

312    of 51:49 (IQR: 55:45, 47:53), which is considerably closer to the human ratio than the mouse ratio.

313    We next used the unpaired $V_H$ sequence data to determine the immunoglobulin heavy-chain gene

314    usage frequencies in Kymice and compare the usage frequency of the IGHV, IGHD, and IGHJ germline

315    genes to those in the human data. We used hierarchical clustering to compare the gene usage profiles

316    of individual Kymice and humans, building phylogenetic trees to show the relationships between the

317    individuals' gene usage profiles. Although Kymice and humans are the same in the numbers of IGHV

318    genes used the frequencies as determined by sequence abundance can differ. The hierarchical

319    clustering of the IGHV genes showed that the Kymice and humans form nearly separate monophyletic

320    clusters except for a single outlier human subject (Figure 2A). Most of the variation in IGHV gene usage

321    is explained by the IGHV gene subgroup usage: clustering by this separates humans and Kymice

322    without the outlier, with Kymice using a lower proportion of IGHV1 and IGHV2 genes compensated

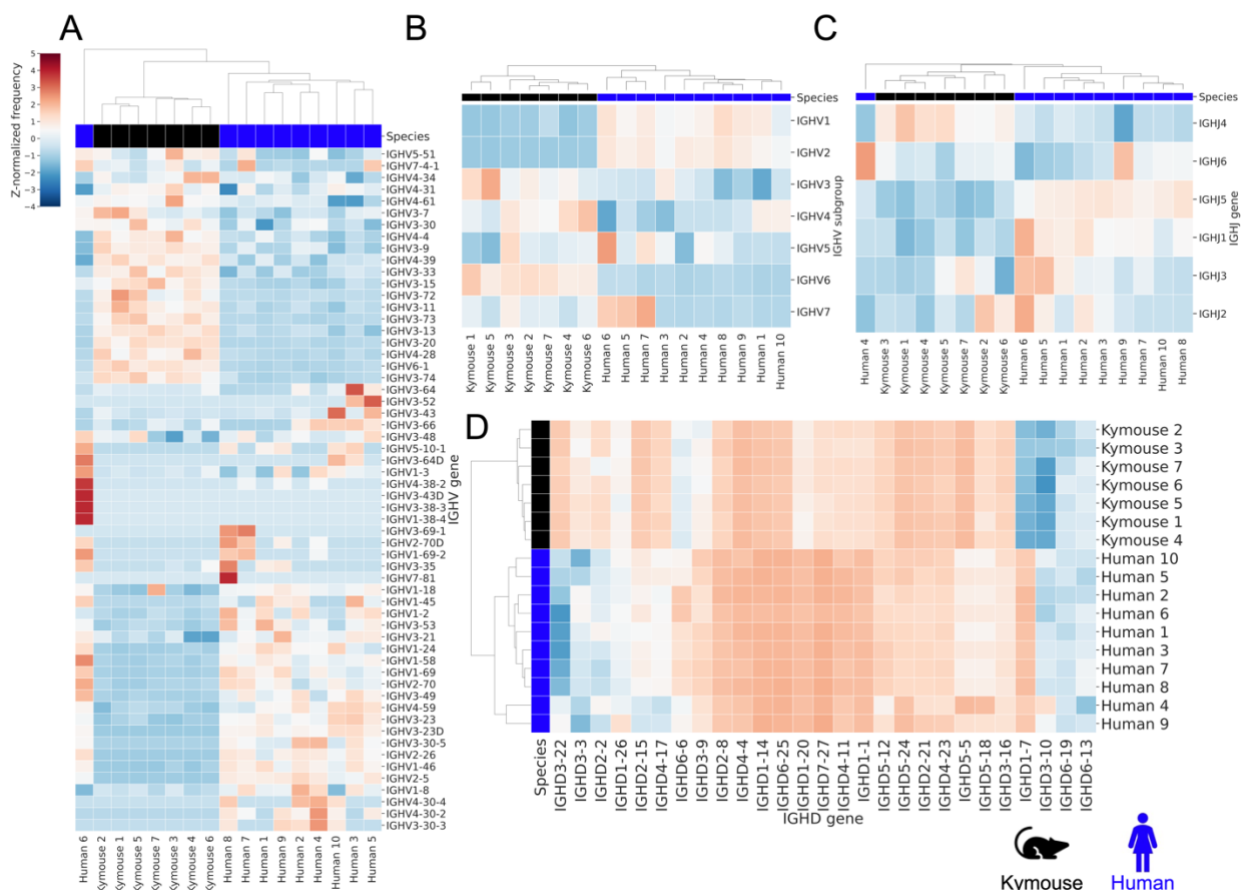323    for by increased IGHV3, IGHV4 and IGHV6 usage (**Figure 2B**).

324

15

325



326 **Figure 2:** *Gene usage clustermaps for A) IGHV genes, B) IGHV subgroups, C) IGHJ genes and D) IGHD*

327 *genes. The IGHV clustermaps show a separation between human (blue) and Kymouse (black)*

328 *repertoires, with lower usage of IGHV1 and IGHV2 in the Kymouse. Downstream, there are also*

329 *differences in usage of IGHJ genes with a preference in the Kymouse repertoires for IGHJ4. The IGHD*

330 *gene usage shows the clearest distinction.*

331

332 The IGHJ gene usage profile is similar: Kymice and eight out of the ten humans form monophyletic

333 clades with two outliers, including the outlier human subject from the IGHV gene usage clustering. On

334 average, the Kymouse uses IGHJ4 more frequently than humans, and IGHJ5 and IGHJ1 less frequently

335 (**Supplementary Figure 1**). Both the IGHV and IGHJ gene usage profiles of naïve Kymouse repertoires

336 are more similar to one another than to any human repertoire. However, Kymice are closer to some

16

337    of the human repertoires than others. This therefore indicates Kymice fall within the range of human

338    IGHV and IGHJ usage frequencies.

339    The hierarchical clustering of the IGHD genes revealed that the Kymice formed a monophyletic sister

340    group distinct to that of the humans (**Figure 2e**). As can be seen from the heatmap, the IGHD germline

341    genes used by the Kymice (*e.g.* IGHD3-22, IGHD2-15) are infrequently used by humans and *vice versa*.

342    We compared the distribution of the CDRH3 lengths in each species' naïve repertoire (**Figure 3A**). We

343    calculated the length of the CDRH3 loop (under the IMGT antibody numbering scheme) from the

344    Kymouse, human and C57BL/6 mouse heavy-chain sequences. This revealed that Kymice have an

345    average CDRH3 length in between that of humans and mice, with a mean CDRH3 length of 14.25. In

346    comparison, the C57BL/6 mouse dataset have a mean length of 12.39 amino acids, and the human

347    dataset have a mean CDRH3 length of 16.56 amino acids.  Kymouse CDRH3 loops are on average 2.36

348    aa shorter than humans (95% CI: 2.26, 2.48; p<0.001), while WT mice CDRH3 loops are on average

349    4.21 aa shorter than humans (95% CI: 4.12, 4.30; p<0.001).

350

351    The length of a CDRH3 loop is determined by four main factors: (i) the choice of IGHV gene; (ii) the

352    choice of IGHD gene; (iii) the choice of IGHJ gene; (iv) the number of nucleotides inserted or deleted

353    in the V-D and D-J junctional regions during B cell development. In order to investigate further why

354    the Kymouse tends to have shorter CDRH3 loops on average than humans, we considered each of

355    these factors in turn. We first looked at whether the differential IGHV gene usage by Kymouse had a

356    significant impact on the length of the CDRH3 loop. For each heavy-chain sequence in the human and

357    Kymouse datasets we determined the number of nucleotides that the IGHV germline gene contributes

358    to the CDRH3 loop. The results showed a statistically significant but small difference of 0.14 nt (95%

359    CI: 0.12, 0.16; p<0.001) between humans and Kymouse. Therefore, it does not appear that the

360    differential choice of IGHV between humans and Kymouse greatly affects the CDRH3 length. We next

361    investigated the effect of the differential usage of IGHD genes between humans and Kymice on the

362    length of the CDRH3 loop. This showed that the human IGHD germline genes used by the Kymouse

363    are, on average, 2.34 nt shorter than humans (95% CI: 2.26, 2.41; p<0.001). We then compared the

364    relative usage of the IGHJ genes of humans and Kymice (**Supplementary Figure 1**). Kymice tended to

365    use IGHJ4 (47.3% in Kymice versus 44.2% in humans) and IGHJ6 (32.8% in Kymice versus 28.6% in

366    humans) while using the other genes (IGHJ1, IGHJ2, IGHJ3, IGHJ5) slightly less frequently, in particular

367    IGHJ5 (11.1% in Kymice versus 15.1% in humans). Estimation statistics revealed that the differential

368    IGHJ gene usage between Kymouse and human resulted in a decrease in the CDRH3 length for

369    Kymouse of 0.85 nt (95% CI: 0.77, 0.94; p<0.001).

370    Finally, we looked at the number of nucleotide insertions in the V-D and D-J junctions in the antibody

371    heavy chains. The results showed that Kymice V-D junctions are on average 3.73 nucleotides shorter

372    than humans (95% CI: 3.68, 3.78 p<0.001), with a mean insertion size of 3.35 nucleotides compared

373    to 7.30 nucleotides for humans (**Supplementary Figure 2A**). Equally, the Kymice D-J junctions are on

374    average 3.68 nucleotides shorter than humans (95% CI: 3.63, 3.73; p<0.001), with a mean insertion

375    size of 2.91 nucleotides compared to 6.77 nucleotides for humans (**Supplementary Figure 2B**). Overall,

376    the number of junctional nucleotides inserted in the Kymouse heavy chain is 7.33 fewer than in

377    humans (95% CI: 7.25, 7.40; p<0.001), with an average of 7.55 in Kymice compared to 14.88 in humans.

378    These results show that the main factors that give rise to the shorter CDRH3 lengths in Kymouse

379    compared to humans are the reduced number of nucleotide insertions in the V-D and D-J junctional

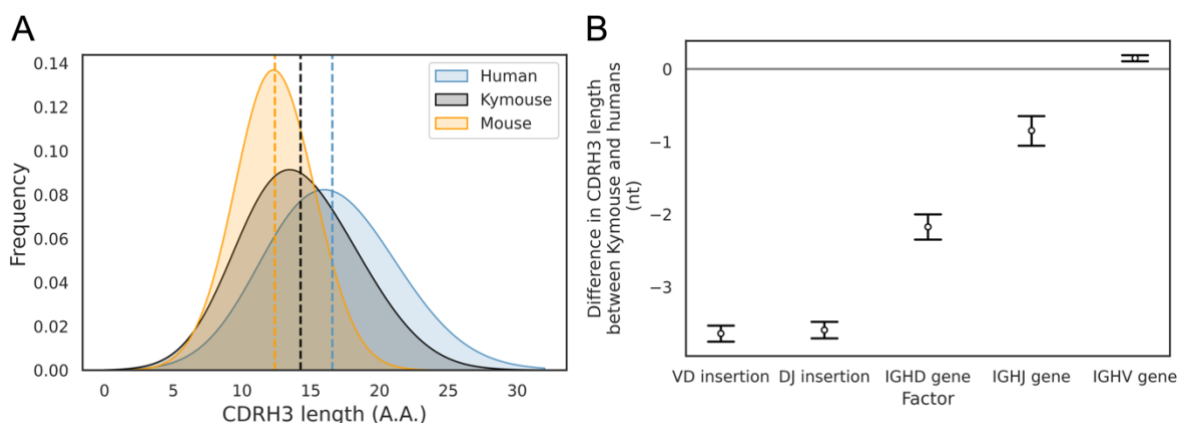380    regions, and the differential usage of IGHD germline genes between the species (**Figure 3B**).



381

382 **Figure 3:** *The CDRH3 length distribution among the mouse, Kymouse and human sequences. The*

383 *Kymouse CDRH3 length distribution is intermediate between the mouse and human distributions (A).*

384 *Plus or minus indicates the standard error of the mean. We used estimation statistics to reveal that*

385 *the major factor leading to this reduction in CDRH3 length despite access to the same germline*

386 *repertoire in Kymouse is the relative lack of VD and DJ insertions (B)(Supplementary Figure 3).*

387 A correlate of shorter CDRH3 length is a reduction in the theoretical CDRH3 diversity obtainable, with

388 each amino acid less in the CDRH3 leading to a further 20-fold reduction in the number of possible

389 CDRH3s. We do not expect the difference in diversity to approach this magnitude because of the

390 CDRH3 sequence patterns created by combinations of germline genes, however a reduction in the

391 non-templated insertion of nucleotides could lead to less diverse repertoires in the Kymouse.

392 To examine the diversity of the repertoires, we considered how many unique CDRH3s or clonotypes

393 can be found in a given repertoire, normalized by the total number of sequences observed. **Figure**

394 **4A** shows that the ratio of the number of unique CDRH3s to the number of sequences is comparable

395 in Kymice and humans, and that both are considerably higher than in C57BL/6 mice. When clustering

396 similar CDRH3s in combination with common IGHV and IGHJ genes (clonotypes), there are

397 comparable unique clonotypes per sequence in mice, humans and the Kymouse (**Figure 4B**). We

398 calculated the Shannon diversity of the CDRH3 and clone distributions, a measure which takes into

399 account both richness and abundance and found comparable diversity across all repertoire types

400 (**Figures 4C and 4D**). The diversity observed depends on CDRH3 length considered, with the mouse

401 and Kymouse repertoires having greater diversity at shorter lengths (peak diversity at lengths 12 and

402 13 respectively) and the human repertoires having greater diversity at lengths greater than 18 amino

403 acids (peak diversity at length 15) for both CDRH3s (**Figure 4E**) and clonotypes (**Figure 4F**). Given

404 that junctional insertions are required to reach lengths of 23 or greater in the human and Kymouse

405 data, this supports the hypothesis that the reduced junctional diversification in naïve Kymouse

406 repertoires limits CDRH3 diversity at the longest lengths, however the greater diversity at shorter

19

407  lengths may compensate. It is also clear that the Kymouse occupies a region in CDRH3 diversity and

408  length between wild type mice and humans.

409  Finally, we considered the CDRH3 and clonotype overlap **(Figures 4G and 4H).** The overlap among

410  CDRH3s is highest between individual mice and individual Kymice. On average 5.1% of CDRH3s were

411  shared between any two Kymice and 3.4% of CDRH3s shared between any two mice compared to

412  0.12% of CDRH3s shared between any two human subjects. Despite the greater CDRH3 diversity

413  reported for individual Kymice, over 25 times more CDRH3s are shared between pairs of Kymice on

414  average than between humans. Clonotype sharing was also higher between individual Kymice

415  (average 7.92% vs 0.84% in humans).

416

417



418  ***Figure 4:*** *Examination of diversity of human, Kymouse and murine repertoires. The top rows pertain*

419  *to exact CDRH3 (amino acids) and the bottom rows to clonotypes (same IGHV, IGHJ and greater than*

420  *90% amino acid identity across length-matched CDRH3s). At the level of CDRH3s, the Kymouse*

421  *repertoires have more unique CDRH3s per sequence sampled (A) are more diversity in their usages*

422  *(C), despite their limited VD and DJ insertion rates. Diversity is reduced relative to human sequences*

423  *at longer CDRH3 lengths which in unmutated repertoires require VD/DJ insertions to reach (E).*

424  *Kymouse repertoires show an opposite pattern in unique sampling rate and diversity when looking at*

20

425    *clonotypes (B and D respectively) but still show reduced diversity versus humans at longer lengths.*

426    *Overlap among CDRH3s (G) and clonotypes (H) between individuals is considerably higher in Kymice*

427    *than humans, and more comparable to mice (G).*

428

429    **The Kymouse repertoires are structurally more human-like than mouse-like**

430    Ultimately, genetic diversity is reflected in the structure of the resulting BCR and secreted antibody

431    protein. Therefore, we compared the structural similarities of the BCRs via structural annotation of

432    the CDRs. We first compared the CDRH1 and CDRH2s, which adopt a limited set of conformations

433    known as canonical classes. For both CDRH1 and CDRH2 the Kymouse repertoires group separately

434    from human repertoires, but group with human repertoires before the mouse repertoires, in their

435    usage of canonical forms (**Supplementary Figure 3**). All canonical forms are strongly predicted by IGHV

436    germline gene subgroup (**Supplementary Table 1**) especially as these sequences are non-mutated.

437    Despite the effects of the differential IGHV usage observed in humans and Kymice, this does not make

438    the Kymouse canonical class usage more similar to mice as the murine repertoires use different

439    canonical classes than either the human or humanised repertoires. Interestingly, six of the nine murine

440    IGHV germlines encode a single canonical form, H1-8-A, suggesting more limited structural diversity

441    in murine CDRH1s.

442    We then performed structural comparisons of the shapes of the human, Kymouse and mouse CDRH3s.

443    CDRH3s do not adopt canonical conformations, so we used two different approaches: firstly, structural

444    annotation which consists of comparison to a CDRH3 structural database and annotation with a

445    structural cluster ID, and secondly full CDRH3 modelling.

446    Of the 1,944 possible structural clusters in the CDRH3 structural database, 1,594 were observed in at

447    least one repertoire. The majority (1,270) of these clusters were observed in all species. There was an

448    observable difference among repertoires in the species origin of the structural clusters observed, i.e.

21

449     each species was biased in the structural space it tended to use. The majority of structural clusters

450     used in human and Kymouse repertoires were of human origin (57.5 and 55.0% respectively), while

451     the majority of structural clusters used in mouse repertoires were of murine origin (64.0%)

452     (**Supplementary Figure 4)**.

453     As described, the CDRH3 lengths in the human dataset were on average 2.36 aa longer than those in

454     the Kymouse, and 4.21 aa longer than those in the C57BL/6 mice (**Figure 3A**). As we considered only

455     CDRH3s between length 4 and 16 in this structural analysis, this difference was reduced to a difference

456     of 0.45 amino acids between human and Kymouse repertoires (CI 0.43, 0.47) and 0.63 between human

457     and mouse repertoires (CI 0.62, 0.65). The average CDRH3 lengths of the human templates in the

458     FREAD database were longer than those in the murine templates (12.6 aa compared to 11.2 aa

459     respectively) (**Supplementary Figure 5**). We checked that the observed differences in structural

460     template usage by humans, Kymice, and C57BL/6 mice were not just reflecting differences in the

461     availability of templates at the difference CDRH3 lengths. Therefore, we stratified the FREAD database

462     by CDRH3 length and ran the datasets against each CDRH3 length separately. If the preferential usage

463     of human-specific templates by humans and humanized mice, and vice-versa for mice, is simply a

464     result of random sampling from the FREAD database, we expect the ratio of the proportion of human

465     templates in human repertoires to the proportion of human templates in the FREAD database to be

466     approximately 1 across all CDRH3 lengths, and below 1 across all CDRH3 lengths for mice. Instead, we

467     see a significant enrichment (at the 1% level) of species-specific templates at multiple CDRH3 lengths.

468     The human repertoires are more structurally variable than are the humanized murine or murine

469     repertoires (**Figures 5A, B and C**). The murine repertoires are the least variable (there is the smallest

470     range of distances between any two given subjects for mice, and on average the smallest average

471     distance between pairs of subjects). All the Kymouse repertoires were structurally closer (Euclidean

472     distance in Z-normalized CDRH3 structural cluster usage) to any given human repertoire than to any

473     C57BL/6 mouse repertoire **(Figures 5A, B and C).** With no correction for sample size or for CDRH3

22

474   length distribution, the humanized murine and human repertoires form a monophyletic cluster that is

475   sister group to the murine repertoires (**Figure 5D**).

476   Labelling cluster members as "murine" or "human-like", the adjusted Rand index of this clustering is

477   1.0 (perfect correspondence). We calculated Rand indices for 100 subsamples with equal numbers of

478   each CDRH3 length to control for any length effect; in all subsamples, the human and humanized

479   repertoires could be clustered separately from the murine repertoires resulting in an adjusted Rand

480   index of 1.0. Further, in all 10 subsamples, at least one human repertoire was more structurally similar

481   to a given humanized murine repertoire than to at least one other human repertoire (between 5 and

482   9 of 10 subjects per subsample).

483   In conclusion, even when adjusting for sample size and differing CDRH3 length distribution, the

484   repertoires of Kymouse are structurally more similar to human repertoires than they are to murine

485   repertoires according to the homology-based structural annotation technique.



486

23

487     *Figure 5: Distance between CDRH3 structural cluster usage profiles is measured with Euclidean*

488     *distance in Z-normalized proportions. This is calculated pairwise between subjects and these distances*

489     *are clustered hierarchically. Figures A through C show these pairwise distances stratified by the type*

490     *of comparison. The leftmost box shows the range in distance between individuals of the same species.*

491     *For the mouse and Kymouse repertoires, this range is smaller than the range in distances for any other*

492     *species, meaning that they cluster monophyletically (D). Human repertoires have less self-similar*

493     *CDRH3 structural cluster usages with ranges overlapping with Kymouse repertoires. In the hierarchical*

494     *clustering solution with these pairwise distances that is shown in D, the human and Kymouse*

495     *repertoires form a monophyletic clade separately from the murine repertoires.*

496

497     The species bias we observed in structural cluster usage, which is assigned via sequence similarity to

498     a structural template, meant that a template-free modelling procedure might lead to different results.

499     We modelled non-singleton clonotype representatives from humans, Kymice and mice using

500     ABlooper, then compared the resultant CDRH3s via $C_\alpha$ RMSD. Analogously to clustering of the CDRH3

501     template database, we performed greedy clustering with a 0.6Å cut-off. This clustered the 43,378

502     models with 41,397 unique CDRH3s into 6,546 clusters. The modelled human CDRH3s were on average

503     longer than Kymouse CDRH3s by 0.99 amino acids (CI: 0.94, 1.05) and longer than mouse CDRH3s by

504     1.23 amino acids (CI: 1.17, 1.29). We observed a slightly different clustering by usage than via the

505     homology approach (**Supplementary Figure 6A**)**,** with monophyly of Kymice, and Kymice and humans,

506     but not of mice. The distribution of intersubject differences was such that the human/Kymouse and

507     human/human comparison distributions are largely overlapping. Indeed, the differences between all

508     distributions were less extreme with overlap between most of the distributions **(Supplementary**

509     **Figure 6B**)**.** The monophyly of Kymice and humans versus mice is observed when subsampling in order

510     to equalize the CDRH3 length distribution, i.e., it is not driven by differences in CDRH3 length. The

511     deep learning-based modelling approach supports the earlier finding that Kymice are more

512  structurally similar to humans than mice, and vice-versa; the extent of this similarity is greater than

513  observed with the homology modelling approach.

514  In conclusion, the sequence differences in the repertoires which were described in the germline and

515  diversity sections do impact upon the repertoire of CDR structures which are observed in the Kymouse.

516  The distribution of canonical forms in CDRH1 and CDRH2 is distinguishable, but usage is human-like.

517  In the CDRH3, both structural annotation and full structure prediction indicate that the naïve Kymouse

518  CDRH3 structural space is human-like, indicating that Kymice repertoires offer comparable structural

519  starting points for the production of antigen-specific antibodies.

520

521  **A state-of-the art humanization tool scores Kymouse sequences as fully human**

522  We next tested whether the Kymouse sequences are considered human by a state-of-the-art

523  humanization tool. We used the random forest classifiers within Hu-mAb to score heavy chain amino

524  acid sequence humanness. 100% of human and Kymouse heavy chain sequences were classified as

525  human, with the maximum humanness score assigned to 99.1% of human sequences and 98.3% of

526  Kymouse sequences. All sequences produced scores in the "Positive (High Score)" category which had

527  minimal anti-drug antibody events reported among therapeutic antibodies. No murine sequences

528  were classified as human.

529

530  # Discussion

531  The phenotypic diversity of B cells in the spleens, lymph nodes and bone marrows of the Kymouse,

532  determined by immunophenotyping panels, showed the main immunologically relevant B cell

533  subpopulations could be identified at appropriate cell frequencies, consistent with the Kymouse being

534  fully competent for B cell development and capable of a complete humoral immune response. As

535    expected, the baseline levels of the immune relevant subsets, i.e. memory B cells and germinal centre

536    B cells in the spleen and lymph nodes as well as the plasma cells in bone marrow were low reflecting

537    the lack of immune exposure beyond commensal and environmental antigens during mouse

538    husbandry.

539    B-cells recognise antigens via the B-cell receptor (BCR) and an individual is able to generate a set of

540    high-affinity BCRs to any given antigen due to the exceptional genetic and structural diversity of B-cell

541    receptor binding sites created by recombination of germline genes, junctional diversification and

542    somatic hypermutation. While somatic hypermutation plays a key role in the development of mature

543    high affinity antibodies, the breadth of an immune response to an antigen is limited at first by the

544    diversity produced by recombination and non-template additions of the non-somatically

545    hypermutated, germline encoded heavy chain immunoglobulin genes and their pairing with similarly

546    rearranged immunoglobulin light chain genes that together comprise the BCRs in the naïve B-cells.

547    Paired $V_H$/$V_L$ sequencing of the naïve Kymouse BCR repertoire revealed a near 50:50 Ig$\kappa$/Ig$\lambda$ ratio

548    which more closely approximates the human Ig$\kappa$/Ig$\lambda$ ratio of 60:40 than the murine Ig$\kappa$/Ig$\lambda$ ratio of

549    95:5. Sequence analysis of immunoglobulin heavy chains showed that naïve Kymouse BCRs have a

550    CDRH3 length distribution that is intermediate between human and mouse repertoires (on average

551    ~14 aa versus ~16 aa for human repertoires and ~12 aa for mice). We compared the IGHV, IGHJ and

552    IGHD gene usage profiles of human and Kymouse repertoires and showed that the IGHV and IGHJ gene

553    usage profiles in Kymice are distinct from human profiles but human-like, while the IGHD gene usage

554    profiles are sufficiently different that Kymouse and human repertoires are distinct. This contrasts with

555    the previous NGS characterization of the OmniRat in which both IGHV and IGHD gene usage was

556    distinct from humans (Joyce, Burton, and Briney 2020). It is an ongoing concern that different germline

557    gene distributions may affect how representative transgenic models are of immune responses to

558    germline-targeting immunogens.

559 We then examined the extent to which the differing germline gene usage distributions contribute to

560 the shorter CDRH3s observed in Kymice. The Kymouse repertoires display a preference for shorter

561 IGHD genes: this was also observed in the NGS characterization of Omni-Rat BCR repertoires

562 suggesting that a preference for shorter IGHD genes may be common across transgenic rodent

563 platforms (Joyce, Burton, and Briney 2020). While the differences in germline gene usage distributions

564 do appear to contribute to the differing CDRH3 length distributions, a greater proportion of the effect

565 is ascribable to fewer nucleotide insertions at both the V(D) and (D)J junctions during junctional

566 diversification in the Kymouse with on average 7.33 nt fewer inserted in Kymouse over the two

567 junctions. This reduced junctional diversification in the Kymouse leads to lower diversity in longer

568 CDRH3s and greater clonotype overlap between individual Kymice than between individual humans.

569 Changes in CDRH3 structural cluster usage have been previously observed along the B cell

570 development axis (Kovaltsuk et al. 2020) and allow comparison of repertoires derived from different

571 species. Despite these genetic differences, modelled structural comparison of the human and

572 humanized repertoires to the murine repertoires, by annotating the $V_H$ sequences with predicted

573 CDRH3 structural template clusters showed murine repertoires use mostly CDRH3 structural clusters

574 that have been identified from murine antibodies, while human and Kymouse repertoires use CDRH3

575 structural clusters identified from a more even distribution of species of which more than 50% were

576 identified from human antibodies. Further, grouping of the exact distribution of CDRH3 structural

577 clusters reveals that Kymouse structural cluster usage is closer to human usage than murine usage,

578 accounting for CDRH3 length differences. When modelling the CDRH3s (as opposed to performing an

579 approximation via structural annotation), part of the "structural distance" between the human, mouse

580 and Kymouse repertoires disappeared suggesting part of the signal observed in homology modelling

581 is due to different sequences with similar predicted shapes. This suggests that despite the observed

582 differences at the sequence level, the CDRH3 structural shapes adopted by the BCRs are within the

583 distribution of observed human shapes.

584    Finally, using the Hu-mAb humanness classifiers, all Kymouse and human sequences are classified as

585    human, meaning that naïve sequences isolated from the Kymouse are predicted to have similar

586    immunogenicity in humans to sequences isolated from humans themselves.

587    In conclusion, although naïve BCR repertoires of the Kymouse have key distinctions from human

588    repertoires at the sequence level they are comparable to the human repertoires in terms of CDRH3

589    structural usage. A number of studies have shown how the Kymouse is able to elicit equivalent

590    antibodies to those found in humans exposed to the same antigen (Sok et al. 2016; Scally et al. 2017;

591    McLeod et al. 2019; Oyen et al. 2020). This suggests that the engagement of BCRs on naïve B cells is

592    authentic, and that the structural templates available for antigen binding are indeed human-like as we

593    show here.

594    **Funding**

598

599

600

601

602

603

604

# References

Abanades, Brennan, Guy Georges, Alexander Bujotzek, and Charlotte M Deane. 2022. 'ABlooper: Fast Accurate Antibody CDR Loop Structure Prediction with Accuracy Estimation'. *Bioinformatics* 38 (7): 1877–80. https://doi.org/10.1093/bioinformatics/btac016.

Allen, Christopher D. C., Takaharu Okada, and Jason G. Cyster. 2007. 'Germinal-Center Organization and Cellular Dynamics'. *Immunity* 27 (2): 190–202. https://doi.org/10.1016/j.immuni.2007.07.009.

Brüggemann, M., H. M. Caskey, C. Teale, H. Waldmann, G. T. Williams, M. A. Surani, and M. S. Neuberger. 1989. 'A Repertoire of Monoclonal Antibodies with Human Heavy Chains from Transgenic Mice'. *Proceedings of the National Academy of Sciences of the United States of America* 86 (17): 6709–13. https://doi.org/10.1073/pnas.86.17.6709.

Brüggemann, Marianne, Michael J. Osborn, Biao Ma, Jasvinder Hayre, Suzanne Avis, Brian Lundstrom, and Roland Buelow. 2015. 'Human Antibody Production in Transgenic Animals'. *Archivum Immunologiae et Therapiae Experimentalis* 63 (2): 101–8. https://doi.org/10.1007/s00005-014-0322-x.

Choi, Yoonjoo, and Charlotte M. Deane. 2010. 'FREAD Revisited: Accurate Loop Structure Prediction Using a Database Search Algorithm'. *Proteins: Structure, Function, and Bioinformatics* 78 (6): 1431–40. https://doi.org/10.1002/prot.22658.

Collins, Andrew M., and Corey T. Watson. 2018. 'Immunoglobulin Light Chain Gene Rearrangements, Receptor Editing and the Development of a Self-Tolerant Antibody Repertoire'. *Frontiers in Immunology* 9: 2249. https://doi.org/10.3389/fimmu.2018.02249.

Corcoran, Lynn M, and David M Tarlinton. 2016. 'Regulation of Germinal Center Responses, Memory B Cells and Plasma Cell Formation—an Update'. *Current Opinion in Immunology*, Lymphocyte development and activation * Tumour immunology, 39 (April): 59–67. https://doi.org/10.1016/j.coi.2015.12.008.

Curtis, Nicholas C., and Jiwon Lee. 2020. 'Beyond Bulk Single-Chain Sequencing: Getting at the Whole Receptor'. *Current Opinion in Systems Biology* 24 (December): 93–99. https://doi.org/10.1016/j.coisb.2020.10.008.

Dunbar, James, and Charlotte M. Deane. 2016. 'ANARCI: Antigen Receptor Numbering and Receptor Classification'. *Bioinformatics* 32 (2): 298–300. https://doi.org/10.1093/bioinformatics/btv552.

636    Dunbar, James, Konrad Krawczyk, Jinwoo Leem, Terry Baker, Angelika Fuchs, Guy Georges, Jiye Shi,

637           and Charlotte M. Deane. 2014. 'SAbDab: The Structural Antibody Database'. *Nucleic Acids*

638           *Research* 42 (D1): D1140–46. https://doi.org/10.1093/nar/gkt1043.

639    Green, Larry L. 2014. 'Transgenic Mouse Strains as Platforms for the Successful Discovery and

640           Development of Human Therapeutic Monoclonal Antibodies'. *Current Drug Discovery*

641           *Technologies* 11 (1): 74–84. https://doi.org/10.2174/15701638113109990038.

642    Greiff, Victor, Enkelejda Miho, Ulrike Menzel, and Sai T. Reddy. 2015. 'Bioinformatic and Statistical

643           Analysis of Adaptive Immune Repertoires'. *Trends in Immunology* 36 (11): 738–49.

644           https://doi.org/10.1016/j.it.2015.09.006.

645    Gupta, Namita T., Jason A. Vander Heiden, Mohamed Uduman, Daniel Gadala-Maria, Gur Yaari, and

646           Steven H. Kleinstein. 2015. 'Change-O: A Toolkit for Analyzing Large-Scale B Cell

647           Immunoglobulin Repertoire Sequencing Data'. *Bioinformatics* 31 (20): 3356–58.

648           https://doi.org/10.1093/bioinformatics/btv359.

649    Jacob, Joshy, Garnett Kelsoe, Klaus Rajewsky, and Ursula Weiss. 1991. 'Intraclonal Generation of

650           Antibody Mutants in Germinal Centres'. *Nature* 354 (6352): 389–92.

651           https://doi.org/10.1038/354389a0.

652    Joyce, Collin, Dennis R. Burton, and Bryan Briney. 2020. 'Comparisons of the Antibody Repertoires of

653           a Humanized Rodent and Humans by High Throughput Sequencing'. *Scientific Reports* 10 (1):

654           1–9. https://doi.org/10.1038/s41598-020-57764-7.

655    Kovaltsuk, Aleksandr, Konrad Krawczyk, Jacob D. Galson, Dominic F. Kelly, Charlotte M. Deane, and

656           Johannes Trück. 2017. 'How B-Cell Receptor Repertoire Sequencing Can Be Enriched with

657           Structural Antibody Data'. *Frontiers in Immunology* 8.

658           https://doi.org/10.3389/fimmu.2017.01753.

659    Kovaltsuk, Aleksandr, Jinwoo Leem, Sebastian Kelm, James Snowden, Charlotte M. Deane, and

660           Konrad Krawczyk. 2018. 'Observed Antibody Space: A Resource for Data Mining Next-

661           Generation Sequencing of Antibody Repertoires'. *The Journal of Immunology* 201 (8): 2502–

662           9. https://doi.org/10.4049/jimmunol.1800708.

663    Kovaltsuk, Aleksandr, Matthew I. J. Raybould, Wing Ki Wong, Claire Marks, Sebastian Kelm, James

664           Snowden, Johannes Trück, and Charlotte M. Deane. 2020. 'Structural Diversity of B-Cell

665           Receptor Repertoires along the B-Cell Differentiation Axis in Humans and Mice'. *PLOS*

666           *Computational Biology* 16 (2): e1007636. https://doi.org/10.1371/journal.pcbi.1007636.

667    Lee, E-Chiang, Qi Liang, Hanif Ali, Luke Bayliss, Alastair Beasley, Tara Bloomfield-Gerdes, Laura

668           Bonoli, et al. 2014. 'Complete Humanization of the Mouse Immunoglobulin Loci Enables

669       Efficient Therapeutic Antibody Discovery'. *Nature Biotechnology* 32 (4): 356–63.

670       https://doi.org/10.1038/nbt.2825.

671  Lefranc, Marie-Paule, Christelle Pommié, Manuel Ruiz, Véronique Giudicelli, Elodie Foulquier, Lisa

672       Truong, Valérie Thouvenin-Contet, and Gérard Lefranc. 2003. 'IMGT Unique Numbering for

673       Immunoglobulin and T Cell Receptor Variable Domains and Ig Superfamily V-like Domains'.

674       *Developmental & Comparative Immunology* 27 (1): 55–77. https://doi.org/10.1016/S0145-

675       305X(02)00039-3.

676  Lemke, A., M. Kraft, K. Roth, R. Riedel, D. Lammerding, and A. E. Hauser. 2016. 'Long-Lived Plasma

677       Cells Are Generated in Mucosal Immune Responses and Contribute to the Bone Marrow

678       Plasma Cell Pool in Mice'. *Mucosal Immunology* 9 (1): 83–97.

679       https://doi.org/10.1038/mi.2015.38.

680  Marks, Claire, and Charlotte M. Deane. 2020. 'How Repertoire Data Are Changing Antibody Science'.

681       *The Journal of Biological Chemistry* 295 (29): 9823–37.

682       https://doi.org/10.1074/jbc.REV120.010181.

683  Marks, Claire, Alissa M Hummer, Mark Chin, and Charlotte M Deane. 2021. 'Humanization of

684       Antibodies Using a Machine Learning Approach on Large-Scale Repertoire Data'.

685       *Bioinformatics*, no. btab434 (June). https://doi.org/10.1093/bioinformatics/btab434.

686  McGuire, K. L., and E. S. Vitetta. 1981. 'Kappa/Lambda Shifts Do Not Occur during Maturation of

687       Murine B Cells.' *The Journal of Immunology* 127 (4): 1670–73.

688  McLeod, Brandon, Kazutoyo Miura, Stephen W. Scally, Alexandre Bosch, Ngan Nguyen, Hanjun Shin,

689       Dongkyoon Kim, et al. 2019. 'Potent Antibody Lineage against Malaria Transmission Elicited

690       by Human Vaccination with Pfs25'. *Nature Communications* 10 (1): 4328.

691       https://doi.org/10.1038/s41467-019-11980-6.

692  Olsen, Tobias H., Fergus Boyles, and Charlotte M. Deane. 2022. 'Observed Antibody Space: A Diverse

693       Database of Cleaned, Annotated, and Translated Unpaired and Paired Antibody Sequences'.

694       *Protein Science* 31 (1): 141–46. https://doi.org/10.1002/pro.4205.

695  Oyen, David, Jonathan L. Torres, Phillip C. Aoto, Yevel Flores-Garcia, Špela Binter, Tossapol

696       Pholcharee, Sean Carroll, et al. 2020. 'Structure and Mechanism of Monoclonal Antibody

697       Binding to the Junctional Epitope of Plasmodium Falciparum Circumsporozoite Protein'.

698       *PLOS Pathogens* 16 (3): e1008373. https://doi.org/10.1371/journal.ppat.1008373.

699  Pantophlet, Ralph, Nino Trattnig, Sasha Murrell, Naiomi Lu, Dennis Chau, Caitlin Rempel, Ian A.

700       Wilson, and Paul Kosma. 2017. 'Bacterially Derived Synthetic Mimetics of Mammalian

701       Oligomannose Prime Antibody Responses That Neutralize HIV Infectivity'. *Nature*

702       *Communications* 8 (1): 1601. https://doi.org/10.1038/s41467-017-01640-y.

703     Raybould, Matthew I. J., Claire Marks, Alan P. Lewis, Jiye Shi, Alexander Bujotzek, Bruck Taddese, and

704          Charlotte M. Deane. 2020. 'Thera-SAbDab: The Therapeutic Structural Antibody Database'.

705          *Nucleic Acids Research* 48 (D1): D383–88. https://doi.org/10.1093/nar/gkz827.

706     Scally, Stephen W., Brandon McLeod, Alexandre Bosch, Kazutoyo Miura, Qi Liang, Sean Carroll, Sini

707          Reponen, et al. 2017. 'Molecular Definition of Multiple Sites of Antibody Inhibition of

708          Malaria Transmission-Blocking Vaccine Antigen Pfs25'. *Nature Communications* 8 (1): 1568.

709          https://doi.org/10.1038/s41467-017-01924-3.

710     Schneider, Constantin, Matthew I J Raybould, and Charlotte M Deane. 2022. 'SAbDab in the Age of

711          Biotherapeutics: Updates Including SAbDab-Nano, the Nanobody Structure Tracker'. *Nucleic*

712          *Acids Research* 50 (D1): D1368–72. https://doi.org/10.1093/nar/gkab1050.

713     Sok, Devin, Bryan Briney, Joseph G. Jardine, Daniel W. Kulp, Sergey Menis, Matthias Pauthner,

714          Andrew Wood, et al. 2016. 'Priming HIV-1 Broadly Neutralizing Antibody Precursors in

715          Human Ig Loci Transgenic Mice'. *Science (New York, N.Y.)* 353 (6307): 1557–60.

716          https://doi.org/10.1126/science.aah3945.

717     Tomayko, Mary M., Natalie C. Steinel, Shannon M. Anderson, and Mark J. Shlomchik. 2010. 'Cutting

718          Edge: Hierarchy of Maturity of Murine Memory B Cell Subsets'. *The Journal of Immunology*

719          185 (12): 7146–50. https://doi.org/10.4049/jimmunol.1002163.

720     Walls, Alexandra C., Brooke Fiala, Alexandra Schäfer, Samuel Wrenn, Minh N. Pham, Michael

721          Murphy, Longping V. Tse, et al. 2020. 'Elicitation of Potent Neutralizing Antibody Responses

722          by Designed Protein Nanoparticle Vaccines for SARS-CoV-2'. *Cell* 183 (5): 1367-1382.e17.

723          https://doi.org/10.1016/j.cell.2020.10.043.

724     Weisel, Florian, and Mark Shlomchik. 2017. 'Memory B Cells of Mice and Humans'. *Annual Review of*

725          *Immunology* 35 (1): 255–84. https://doi.org/10.1146/annurev-immunol-041015-055531.

726     Wong, Wing Ki, Guy Georges, Francesca Ros, Sebastian Kelm, Alan P. Lewis, Bruck Taddese, Jinwoo

727          Leem, and Charlotte M. Deane. 2019. 'SCALOP: Sequence-Based Antibody Canonical Loop

728          Structure Annotation'. *Bioinformatics* 35 (10): 1774–76.

729          https://doi.org/10.1093/bioinformatics/bty877.

730     Xu, John L, and Mark M Davis. 2000. 'Diversity in the CDR3 Region of VH Is Sufficient for Most

731          Antibody Specificities'. *Immunity* 13 (1): 37–45. https://doi.org/10.1016/S1074-
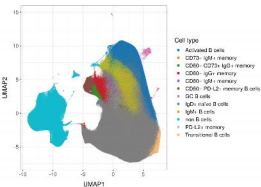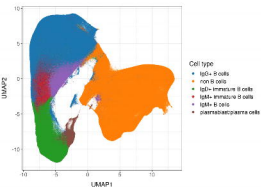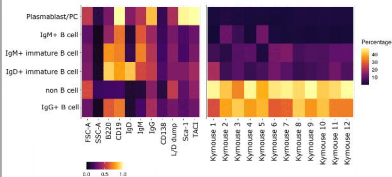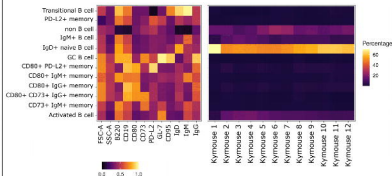
732          7613(00)00006-6.

733     Yaari, Gur, and Steven H. Kleinstein. 2015. 'Practical Guidelines for B-Cell Receptor Repertoire

734          Sequencing Analysis'. *Genome Medicine* 7 (1): 121. https://doi.org/10.1186/s13073-015-
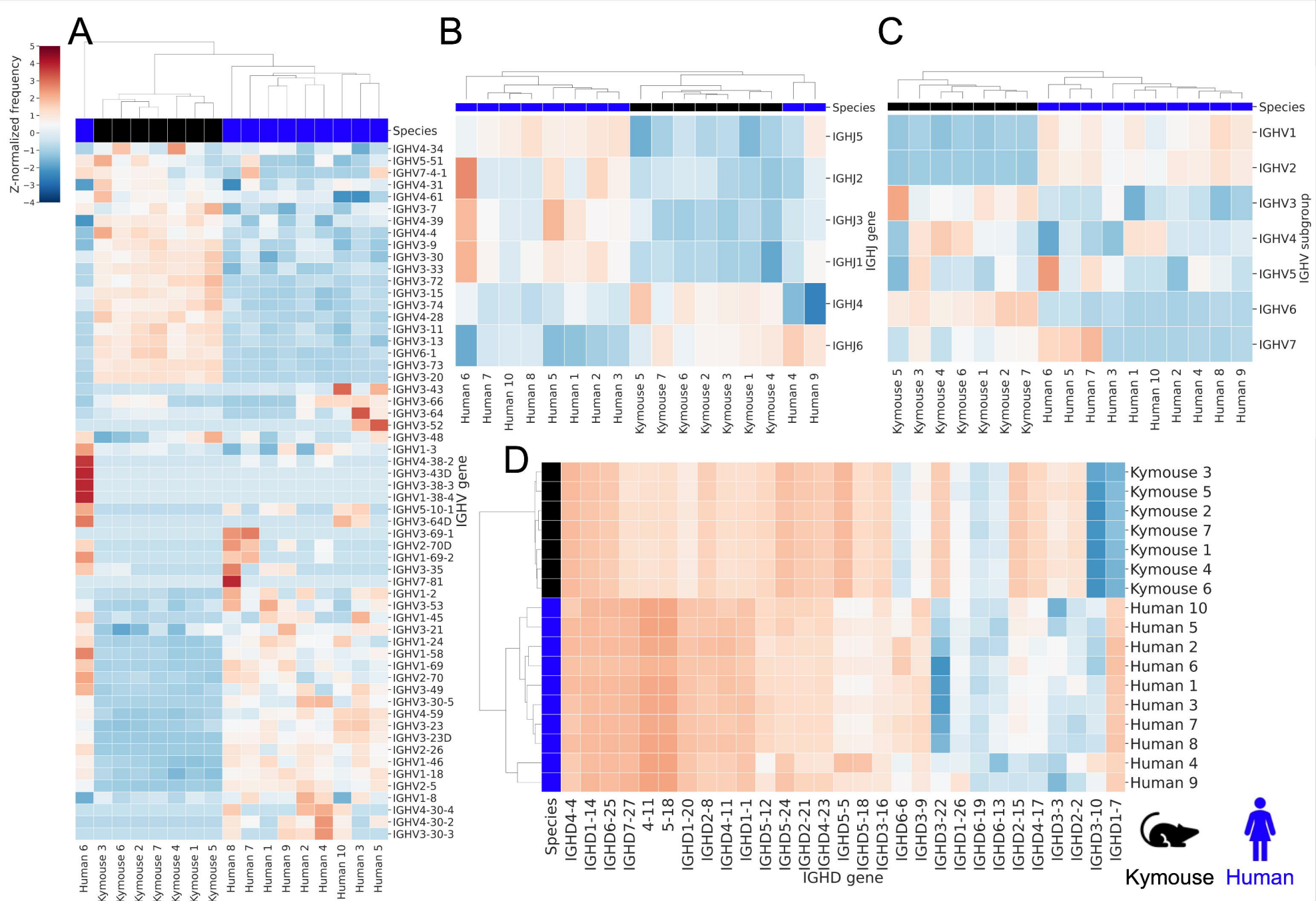
735          0243-2.

736    Ye, Jian, Ning Ma, Thomas L. Madden, and James M. Ostell. 2013. 'IgBLAST: An Immunoglobulin

737        Variable Domain Sequence Analysis Tool'. *Nucleic Acids Research* 41 (W1): W34–40.
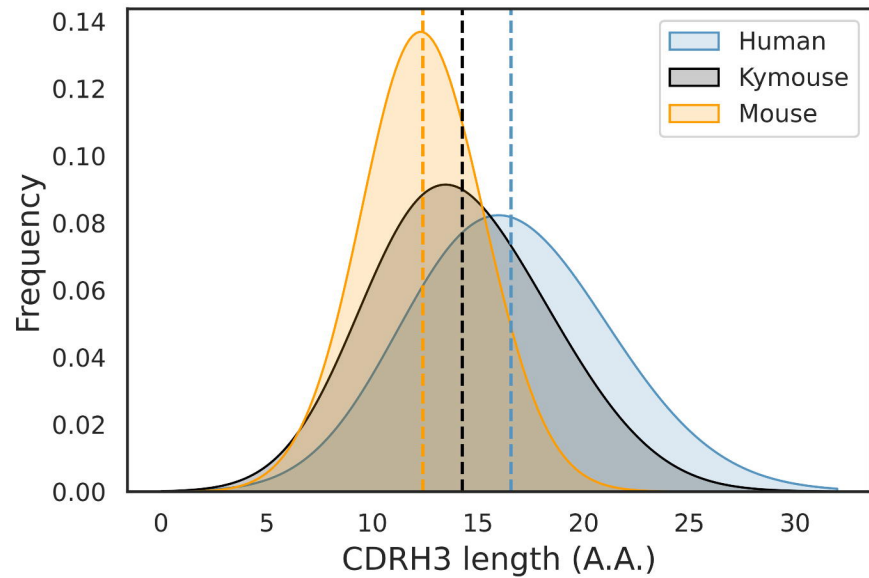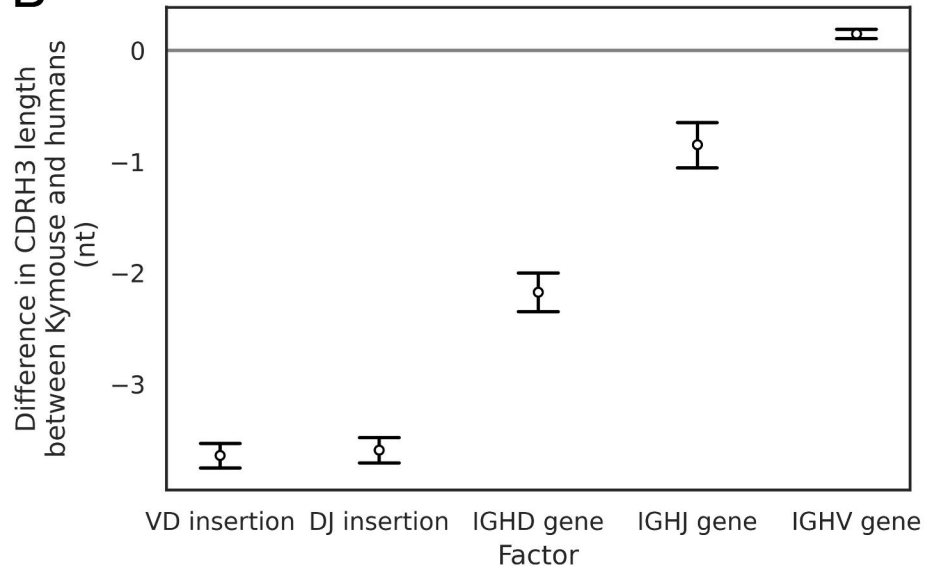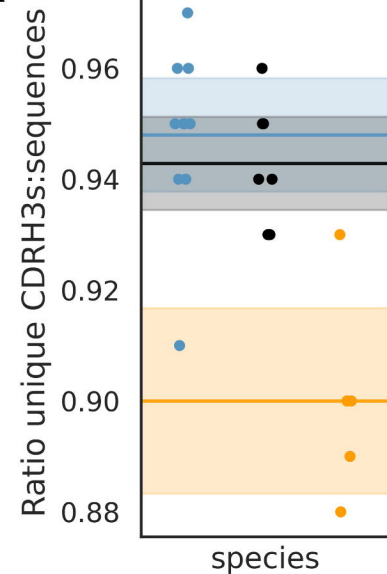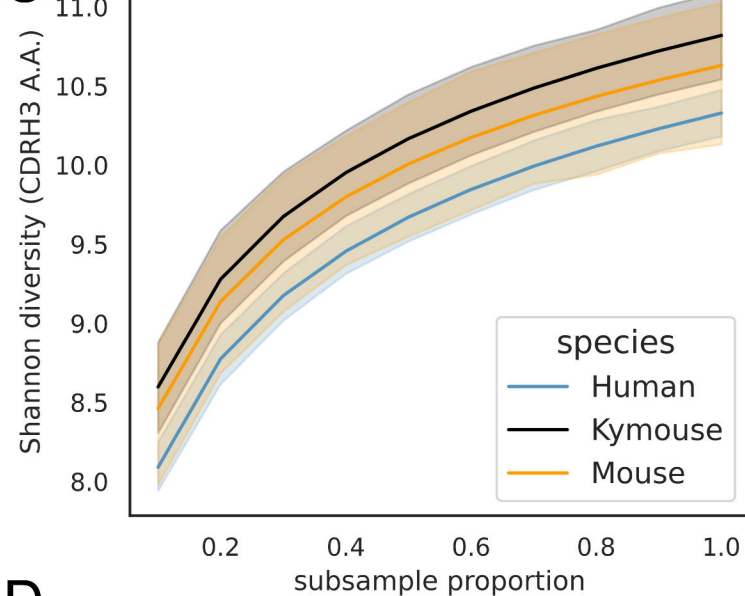
738        https://doi.org/10.1093/nar/gkt382.

739


740


741

**A** Comparisons with murine repertoires

**B** Comparisons with Kymouse repertoires

**C** Comparisons with human repertoires

**D** normalised frequency

species

CDRH3 structural cluster