1

2   DR. ZHUI  TU (Orcid ID : 0000-0003-3187-0132)

3

4

6

7

8

9   **Landscape of Variable Domain of Heavy-chain-only Antibody**

10   **Repertoire from Alpaca**

11

12   Zhui Tu [1,2,3,4], Xiaoqiang Huang [2], Jinheng Fu [1, 5], Na Hu [1, 4, 6], Wei Zheng [2], Yanping Li [1, 4,]

13   [5,]**Error! Bookmark not defined.**, Yang Zhang [2, 3,] **Error! Bookmark not defined.**

14   [1] *State Key Laboratory of Food Science and Technology, Nanchang University, Nanchang,*

15   *China,* [2] *Department of Computational Medicine and Bioinformatics, University of Michigan,*

16   *Ann Arbor, MI, USA,* [3] *Department of Biological Chemistry, University of Michigan, Ann*

17   *Arbor, MI, USA,* [4] *Jiangxi Province Key Laboratory of Modern Analytical Science, Nanchang*

18   *University, Nanchang, China,* [5] *Jiangxi-OAI Joint Research Institution, Nanchang University,*

19   *Nanchang, China, and* [6] *Maternal and Child Medical Research Institute, Shenzhen Maternity*

20   *and Child Healthcare Hospital, Southern Medical University, Shenzhen, China*

21

22   **Running title:** *Tu Z et al / Analysis of Immune Repertoire in Alpaca*

23

24   Abbreviations: HCAbs,  heavy-chain-only antibodies; VHHs, the variable regions of heavy

25   chain of HCAbs; HTS, high-throughput sequencing; CDR, complementary determining

26   region; GSSPs, germline specific scoring profiles; AAs, amino acids; SR, substitution rate;

27   ASR, average substitution rate; SHM, somatic hypermutation; PBMCs, peripheral blood

28   mononuclear cells; PCR, polymerase chain reaction; MSAs, multiple sequence alignments

29

30    *Correspondence: Yanping Li, Jiangxi-OAI Joint Research Institute, Nanchang University,

31    235 East Nanjing Road, Nanchang, Jiangxi, 330047, China. Email: liyanping@ncu.edu.cn (Li

32    Y); and Yang Zhang, Department of Computational Medicine and Bioinformatics, University

33    of Michigan, 100 Washtenaw Avenue, Ann Arbor, MI 48109-2218, USA. Email:

34    zhng@umich.edu (Zhang Y). Senior author: Yang Zhang

## Summary

Heavy-chain-only antibodies (HCAbs), which are devoid of light chains, have been found naturally occurring in various species including camelids and cartilaginous fish. Due to their high thermostability, refoldability, and capacity for cell permeation, the variable regions of heavy chain of HCAbs (VHHs) have been widely used in diagnosis, bio-image, food safety, and therapeutics. Most immunogenetic and functional studies of HCAbs are based on case studies or a limited number of low throughput sequencing data. A complete picture derived from more abundant high-throughput sequencing (HTS) data can help us gain deeper insights. Thus, we cloned and sequenced the full-length coding region of VHHs in Alpaca (*Vicugna pacos*) via HTS in this study. A new pipeline was developed to conduct an in-depth analysis of the HCAbs repertoires. Various critical features, including the length distribution of complementary determining region 3 (CDR3), V(D)J usage, VJ pairing, germline specific mutation rate, and germline specific scoring profiles (GSSPs), were systematically characterized. The quantitative data show that V(D)J usage and VHHs recombination are highly biased. Interestingly, we found that the average CDR3 length of classical VHHs is longer than that of non-classical ones, whereas the mutation rates are similar in both kinds of VHHs. Finally, GSSPs were built to quantitatively describe and compare sequences that originate from each VJ pairs. Overall, this study presents a comprehensive landscape of the HCAbs repertoire, which can provide useful guidance for the modeling of somatic hypermutation and the design of novel functional VHHs or VHH repertoires via evolutionary profiles.

**Keywords:** High-throughput sequencing; Nanobody; Immune repertoire; Antibody diversity; Protein design

## Introduction

The antigen binding domain of functional heavy-chain-only antibodies (HCAbs) discovered in camelids and sharks is composed of a single variable domain.[1, 2] The variable regions of heavy chain of HCAbs (VHHs), also known as nanobodies, have attracted growing interest in various applications, as they are more soluble and stable than canonical antibodies (VHs).[3-6] In camels, the ratio of HCAbs to total IgGs can reach more than 80%, which indicates that HCAbs play a significant role in immune protection.[7] However, it is obvious that the diversity of HCAbs is dramatically lower than that of canonical antibodies due to the lack of VH-VL

67  combinational diversification. This raises a question of how HCAbs can compete with
68  canonical antibodies. Several hypotheses and observations have been proposed over the past
69  decades to address the problem of diversity reduction inherent to HCAbs. One hypothesis is
70  that the complementary determining region 3 (CDR3) of VHHs contains longer loops than
71  canonical antibody VHs (18 amino acids versus 13 amino acids), which helps compensate for
72  the lack of diversity.[8] Evidently, longer CDR3 length increases the paratope size, as well as
73  the three-dimensional structural diversity and contact surface area with antigens.[9] Another
74  explanation, inferred from a structural study that compared two independently generated anti-
75  lysozyme nanobodies, is that the *in vivo* maturation and selection systems are strong enough
76  to compensate for the decrease in the VHHs primary repertoire.[10]

77      High-throughput sequencing (HTS) technology enables scientists to evaluate millions of
78  sequences in parallel, resulting in the collection of more complete and comprehensive
79  information for target samples. This capability makes HTS suitable for the characterization of
80  immune repertoires that are highly plastic and diverse. Although HTS is now routinely
81  applied in the studies of human adaptive immunity,[11] vaccine development,[12] and diagnostic
82  research,[13] only a few studies were tried on VHHs. Fridy *et al.* developed a pipeline
83  combining HTS and proteomics to identify specific VHHs.[14] Similarly, Turner *et al.*
84  demonstrated that HTS can be used as a complementary tool for phage-display bio-panning to
85  rapidly obtain additional clones from an immune VHH library.[15] For the first time, Li *et al.*
86  compared the repertoires of classical antibodies and HCAbs of *Bactrian* camels, with analysis
87  data including CDR3 length distribution, mutation rate, characteristic amino acids, the
88  distribution of cysteine codons, and the non-classical VHHs.[8] Nevertheless, the features of
89  HCAbs, such as the germline usage and mutation preferences, still remain unknown. Like
90  classical immunoglobulin (Ig) heavy-chains, VHHs are encoded by recombined V(D)J genes
91  that are formed from sets of Variable (V), Diversity (D), and Joining (J) genes (IGHV, IGHD,
92  IGHJ) on the genome. An in-depth analysis of the origination and mutation profiles of VHHs
93  would help us to better understand the diversity of the HCAb repertoire, as well as the
94  diversity compensation. Furthermore, appropriate interpretation of the information is
95  important to guide the design of novel functional VHHs.[16, 17]

96      This study is mainly focused on the HCAb repertoire. First, the coding sequences of
97  VHHs from long-hinge HCAbs (IgG$_2$) and short-hinge HCAbs (IgG$_3$) are amplified from the
98  non-immunized and the antigen immunized antibody repertoires of *Vicugna pacos*, where
99  full-length coding sequences of VHHs are obtained by an Illumina MiSeq System (2×300)

100 under the paired-end module. Next, a new pipeline combined with multiple software tools is
101 developed to characterize the diversity and evolutionary features of the VHHs, including
102 CDR3 length distribution, V(D)J usage, VJ pairing, DJ pairing, germline specific mutation
103 rate, and germline specific scoring profiles (GSSPs) (Fig. 1). Considering that the diversity of
104 antibody repertoires is position, chain, and species-dependent,[18-20] comparative studies are
105 also made on amino acid sequences derived from different germline genes.

## Materials and methods

107 *RNA extraction and reverse transcription*

108 Peripheral blood mononuclear cells (PBMCs) were separated from peripheral blood by Ficoll-
109 1.077 (Sangon, Shanghai, China) gradient centrifugation, separately. Three naïve blood
110 samples were collected from three non-immunized healthy male Alpaca (*Vicugna pacos*). To
111 collect immunized blood samples, one donor was immunized by subcutaneous, lower-back
112 injections every two weeks. Samples of fresh blood were collected 1 week after the fifth and
113 seventh immunization. For each blood sample, RNA was purified from approximately $2 \times 10^7$
114 PBMCs using RNAprep Kit (Tiangen, Beijing, China), following the manufacturer's
115 instruction. First-strand complementary DNA (cDNA) was synthesized with random hexamer
116 primers using PrimeScript$^{TM}$ RT-PCR Kit (TAKARA, Dalian, China), and then stored at -
117 80 °C.[21]

118 *Library construction and Illumina sequencing*

119 The VHH coding region was amplified from cDNA by a nested polymerase chain reaction
120 (PCR) as described before.[22, 23] In brief, the variable region was first amplified by primers
121 AlpVh-LD and AlpVHH, which anneal to the conserved region of the leading sequence and
122 CH2 region, respectively. Next, the PCR products were diluted as a template for the second
123 round of PCR, which employed primer pairs AlpVHH-F/AlpVHH-R1 and AlpVHH-
124 F/AlpVHH-R2 to amplify coding sequences of short and long hinge heavy chain antibodies,
125 respectively. The PCR products that encoded VHHs (~450 bp) were purified using TAKARA
126 gel extraction kits (Dalian, China), and then subjected to Next-Generation Sequencing by the
127 Beijing Genomics Institute (BGI) sequencing center. Sequences were generated with a MiSeq
128 System using a 2×300 paired-end module.

129 *Basic data processing*

130 Adapter sequences were first checked and removed from the reads. Then, the reads that bases
131 of "N" were greater than 10% or have > 50% bases with quality values ≤ 5 were discarded,

132     resulting in 14.13×2 million paired-end reads. The pairwise reads were joined using the fastq-

133     join tool (version 1.3.1).[24] The main parameters were the maximum difference percentage

134     (8%) and the minimum base overlap (6 bp). Phylogenetic trees of V germline genes were built

135     using MEGA version X.[25]

136     *V(D)J assignment and numbering*

137     The V(D)J germline gene sequences were obtained from the international ImMunoGeneTics

138     information system (IMGT) antibody repertoire database.[26] The out-frame bases in the 3' end

139     of J gene sequences were manually deleted, where the elaborated germline gene sequences

140     were used to build an IgBLAST database. The resulting 5.85 million joined sequences were

141     subjected to IgBLAST1.8.0 with default parameters.[27] The origin of each sequence, either

142     from long-hinge or short-hinge IgGs, was identified by BLAST 2.7.1+ according to the E-

143     value and sequence identity of the alignments.[28] The V(D)J germline genes on the top of the

144     resulting list of IgBLAST were assigned to the sequences. An in-house Python script was

145     used to analyze VHHs length distribution, CDR3 length distribution, V/D/J germline gene

146     usage, VJ pairing, DJ pairing, and amino acid substitution. The IMGT numbering system was

147     adopted for the coding sequences of VHHs.

148     *Construction and comparison of GSSPs*

149     The sequences were translated and aligned to all alleles of the gene. Sequences with more

150     than one stop codon or no amino acid substitution were discarded. Sequences belonging to the

151     same VJ germline gene were parsed from IgBLAST output to build multiple sequence

152     alignments (MSAs); redundant amino acid sequences were removed in each MSAs. To

153     improve the accuracy of sequence alignment, the V and J segments of each amino acid

154     sequence in an MSA were re-aligned with corresponding IMGT numbered germline

155     sequences     using     an     in-house     NW-align     program     (Y.     Zhang,

156     https://zhanglab.ccmb.med.umich.edu/NW-align/) before GSSP construction. The GSSPs

157     were built and compared as described in previous work.[20] In brief, the MSAs whose number

158     of sequences was greater than a given threshold (e.g. 100, 500, and 1000) were used to build

159     GSSPs. were used to build GSSPs, respectively. A divergence matrix between GSSPs was

160     calculated; each element in the matrix was the Jensen-Shannon divergence calculated between

161     each pair of sequences from the MSAs. The R function *cmdscale* was used for

162     multidimensional scaling and generating coordinates for plotting. The logo plot of MSA was

163     drawn using a stand-alone version of WebLogo 3.6. [29]

164     *Calculation of substitution frequencies for the 20 AAs*

165  The GSSPs were used to calculate the substitution frequencies. The substitution rate (SR)

166  from each GSSP is calculated by

167
$$SR = \frac{\Sigma_{i=1}^{N} f_i}{L \times N} \times 100\% \qquad 1)$$

168  where $f_i$ is the mutation frequency of sequence $i$ to the corresponding germline genes. $L$ is the

169  length of the GSSP, and $N$ is the total sequences in the MSA. The average substitution rate

170  (ASR) for the 20 AAs of a GSSP is calculated by

171
$$ASR_{(a,b)} = \frac{\Sigma_{i=1}^{L} fi_{(a,b)}}{f_{(a)} \times N} \times 100\% \qquad 2)$$

172  where $ASR_{(a, b)}$ is the average substitution rate of amino acid $a$ in germline gene substituted by

173  observed amino acid $b$ in MSA, $fi_{(a, b)}$ is the frequency of amino acid $a$ in germline gene

174  substituted by amino acid $b$ at the position $i$ of an MSA, $f_{(a)}$ is the frequency of amino acid $a$

175  in germline sequence, $L$ is the length of the MSA, and $N$ is the total sequences in the MSA.

176  *Statistical analysis*

177  To investigate the likelihood of pairing preference between germline segments, we used an *in*

178  *silico* simulation protocol as described in a previous study.[30] Briefly, in each simulation, an

179  equal number of real data sequences were constructed using the same individual frequencies

180  of V, D, and J segments observed in the real data. After 2,000 simulation steps, the DJ and VJ

181  pairing that appeared in each simulation were counted. The relative deviation (RD) of

182  minimum, maximum, and real frequencies of each kind of pairing were calculated by

183
$$RD = \frac{x - \bar{x}}{\bar{x}} \times 100\% \qquad 3)$$

184  where $x$ is the minimum or maximum frequencies of simulation, or frequencies of real

185  sequence data, and $\bar{x}$ is the average frequency of each pairing in the 2,000 simulation steps.

186      We used the function *spearmanr* in the Python module *scipy* to calculate the Spearman's

187  Rank Correlation Coefficient to evaluate the statistical dependence of the germline usage, VJ

188  and DJ pairings, and the substitution preference between samples.

189  **Results**

190  *Sequence data filtration and formation*

191  A summary of the sequencing datasets processed in this study is shown in Table 1. The MiSeq

192  sequencing of the non-immune and antigen-experienced HCAb repertoires yielded a total of

193  38.25×2 million reads. Since the sample Naïve-1 generated the most sequencing reads

194  (14.13×2 million reads), it was used to build and test the pipeline. A number of 2,550,856

195  unique DNA sequences were subjected to IgBLAST to identify the germline gene origination

196 of each sequence, after the redundant DNA sequences of the joined paired-end reads were
197 removed. Both V and J germline genes are found in more than 97% of the non-redundant
198 DNA sequences. Following these filtrations, a total of 2,490,298 unique DNA sequences with
199 VJ assignment hits were used to determine the coding sequence (CDS) distribution, V(D)J
200 usage, VJ pairing, and DJ pairing. Briefly, the CDS length distribution centers around 375 bp
201 and follows an approximately normal distribution, where the maximum CDS length is 438 bp
202 in the dataset (Fig. S1 in Supplementary Material). A number of 1,973,186 unique amino acid
203 sequences deduced from this dataset were used to construct multiple sequence alignments
204 (MSAs), to analyze CDR3 length distribution, and to calculate substitution rates and construct
205 GSSPs. VHHs from long-hinge and short-hinge HCAbs were identified and analyzed for
206 comparison.

207 *Germline gene usage*

208 Studies of canonical antibody repertoires have demonstrated that specific V, D, and J
209 germline genes have very different frequencies in humans and mice.[30-33] Meanwhile, HCAbs
210 and canonical IgGs in Alpaca (*Vicugna pacos*) genome have been shown to originate from the
211 same IgH locus, which is composed of 88 V genes (including 4 pseudogenes), 8 D genes, and
212 7 J genes.[34] Here, we utilized the tool IgBLAST to determine the origination of V, D, and J of
213 each clone. The 84 functional V genes, 8 D genes, and 7 J genes were employed to create a
214 reference database for IgBLAST. The IgBLAST results showed that the V, D, and J segment
215 usages have strong preferences for specific germline genes (Fig. 2).

216    The V segments of all clones were generated from the subgroups of IgHV3. The V
217 segments IGHV3S65*01, IGHV3S3-3*01, and IGHV3S53*01 are used by more than 10% of
218 all the clones, while the top 11 V germline genes are used by more than 95% of all the clones
219 (Fig. 2A). All the 17 V germline genes, which contain at least two framework region 2 (FR2)
220 hallmark residues, F37, E44, R45, and G47 in the Kabat numbering system,[35] are in a sub-
221 cluster of IgHV3 (Fig. S2 in Supplementary Material). Germline genes from this sub-cluster
222 contribute more than 85% of V gene usage (Fig. 2A). These hallmark residues are considered
223 to be important for the solubility and stability of VHHs, as well as the VH-VL association of
224 conventional VHs. A novel promiscuous class of VHHs that do not have any FR2 imprints
225 was reported in Sanger sequencing studies.[36, 37] It is now clear that sequences that lack FR2
226 imprints are generated from other V germline genes, in which IGHV3S39*01,
227 IGHV3S41*01, IGHV3S25*01, IGHV3-1*01, IGHV3S9*01, and IGHV3S1*01 constitute

228  the top 6 contributors. These hallmark-free V segments are responsible for about 10% of V
229  gene usage in the dataset.

230  The usage of D segments was relatively evenly distributed across the germline genes,
231  where six out of eight D germline genes have above 10% usage (Fig. 2B). Similar to the V
232  gene usage, the J germline gene usage was also highly biased (Fig. 2C). For instance, the
233  germline gene IGHJ4*01 was used by two-thirds of the J segments. Since only a few
234  sequences were assigned to IGHJ5*01 and IGHJ1*01 (0.15% and 0.09%, respectively), we
235  manually checked the DNA and corresponding amino acid sequences. The IGHJ5*01 hits of J
236  segments were correctly assigned by IgBLAST. However, due to the defects in the 3'
237  sequences, all the IGHJ1*01 assignments were false positives, indicating that VHHs never
238  use IGHJ1*01. These sequences were therefore discarded in the subsequent analyses.

239  *V(D)J recombination preferences*

240  VJ pairing data showed that more than 90% of the VJ pairs are composed of genes from the
241  top 21 most used VJ germline gene combinations (Fig. 3A), indicating that VJ pairing is
242  biased. Theoretically, the combination of VJ pairing should be much greater than 21, even
243  though the V and J usage are highly biased toward specific germline genes. To evaluate
244  whether V(D)J pairing exhibits bias, simulated antibody repertoires were employed to test
245  statistical preference. As the V(D)J recombination occurs in two steps to assemble a complete
246  variable region *in vivo*, we firstly analyzed the DJ pairing, and then the VJ pairing. Although
247  most relative deviation of the real data was less than 100%, DJ pairing showed a preference
248  (Fig. 3B). The VJ pairing results indicated a stronger bias, as the relative deviations of the 6
249  types of VJ combination were more than 200% (Fig. 3C). Notably, all the highly biased VJ
250  pairing were from FR2 hallmark free V germline genes, which were IGHV3-1*01,
251  IGHV3S1*01, IGHV3S25*01, and IGHV3S39*01.

252  *CDR3 length and distribution*

253  The CDR3 length of VHHs from the HTS data mainly ranged from 4 to 34 amino acids
254  (AAs), according to the IMGT numbering system (Fig. 4). The overall average length of
255  HCAbs CDR3 is 18 AAs, consistent with previous studies.[8] We found that the shortest and
256  longest CDR3 lengths were 2 and 39 AAs, respectively, although they were quite rare.
257  Interestingly, VHHs derived from various germline genes showed different CDR3 length
258  distributions (Table S1 in Supplementary Material), indicating a bias of insertion during the
259  process of *in vivo* V(D)J recombination. Hence, we further compared the sequences derived
260  from the top 11 V germline genes (Fig. 4). Notably, the results showed that the average CDR3

261 length of clones derived from hallmark germline genes is longer than that of hallmark free
262 germline genes, except IGV3S9*01.

263 *Substitution and insertion analysis*

264 Since CDR3 contains random insertions and is highly diverse, only the VJ paired segments
265 were used for substitution analysis. The substitution rate (SR), which represents the mutation
266 strength of a VJ pair lineage, ranged from 12% to 22% (Table S2 in Supplementary Material).
267 To analyze the substitution preference of each amino acid, we calculated the average
268 substitution rate (ASR) of the VJ pairing that is comprised of more than 1000 lineages, and
269 then overall ASR. The results demonstrated that partial substitutions tend to be biased and
270 most types of mutations are rare (Fig. 5). As to the overall ASR, 79 out of 441 substitution
271 types are higher than 1%. Insertion of glycine (20.78%) and alanine (12.18%) are preferred at
272 the tip of CDR1 and CDR2 loop. Meanwhile, we found that each germline VJ pair showed
273 various substitution patterns. Therefore, we further calculated the germline specific
274 substitution profiles (GSSPs) to quantify and compare the diversity clustered by each
275 germline VJ pair.

276 *Construction and comparison of VHH profiles*

277 A GSSP which captures the frequency at which each amino acid appears at every position in
278 an MSA is an N by L matrix, where N is the number of residue types and L is the alignment
279 length. The weighted average of the Jensen-Shannon divergence between GSSPs was
280 calculated to quantitatively compare different profiles and then visualized using
281 multidimensional scaling. In order to test the robustness of this quantification method for
282 GSSPs, we calculated Jensen-Shannon divergence of VJ pairing types that have more than
283 100, 500, or 1000 lineages, respectively. The results confirmed that lineages from common V
284 genes tend to be clustered, no matter what cutoff values were used (Fig. 6). Moreover,
285 plotting Jensen-Shannon divergence of the top 11 V germline family showed that some
286 classes are close to each other, indicating that mutation patterns are similar between clustered
287 families (Fig. 6 B, D, and F).

288 *Comparison of long-hinge and short-hinge HACbs*

289 Specific primers were designed to amplify $IgG_2$ and $IgG_3$, which enabled the identification of
290 each clone type. The IgG specific primer sequences were found in 5,674,954 sequences (97%
291 of all unique sequences). Comparison of $IgG_2$ and $IgG_3$ showed that the ranks of J gene usage
292 are the same, but ranks of V and D usage are different, indicating a different preference of V
293 and D segments (Fig. 7). Notably, a bias of gene rearrangement was observed for these two

294 types of HCAbs. The top 5 hallmark V germline genes contribute 90.91% of long-hinge

295 (IgG$_2$) clones, but only 75.30% of short-hinge (IgG$_3$).

296 *Comparison of VHHs from different donors*

297 To test the robustness of the pipeline, the HTS sequences from four other peripheral blood

298 samples (Naïve-2, Naïve-3, Immune-1, and Immune-2), which were collected from the non-

299 immunized and immunized donors (Table 1), were processed following the same pipeline

300 respectively. The V and J germline usage, VJ pairing, DJ pairing, and the substitution

301 preference are highly correlated between the five samples (Table S3 in Supplementary

302 Material). Interestingly, the correlations of the D germline usage are low between samples

303 (Table S3 in Supplementary Material), especially between the naïve and the immunized

304 samples (Spearman rank correlation coefficients: Immune-1 and Naïve-1, $\rho = 0.683$, P =

305 0.042; Immune-1 and Naïve-3, $\rho = 0.467$, P = 0.205).

306 **Discussion**

307 HCAbs naturally occur in various species such as camelids (e.g., camels and llamas) and

308 cartilaginous fish (e.g., sharks).[38] This remarkable evolutionary convergence implies the

309 advantages of functional HCAbs. Thus, systematically investigation of HCAb repertoires is

310 important to reveal the mystery of evolutionary conservation as well as to understand the

311 compensation for the lack of diversity in HCAbs. In this study, we developed a novel pipeline

312 to analyze the full coding sequences of variable domains. In order to automatically process

313 data and maintain its reliability, we tried to avoid using arbitrary filters in the workflow when

314 possible and checked the intermediate results from each step. Since MSAs are crucial for

315 calculating substitution rates and building GSSPs, a classic NW-align algorithm was

316 employed to re-align MSAs from IgBLAST. In order to mitigate effects from noise in the

317 data, we set 1000 as the minimum number of lineages to calculate ASR. The pipeline can be

318 easily extended to analyze the HTS data of antibody repertoires from other species.

319 V(D)J recombination is one of the mechanisms of antibody diversity. Previous study

320 confirmed that the V germline genes of HCAbs and conventional IgGs were located in the

321 same IgH locus on the genome.[34] The hallmark residues in FR2 regions have been known as a

322 characteristic of VHHs. A previous study reported the presence of novel hallmark-free

323 variable domains that can be rearranged to both camelid classical antibodies and HCAbs.[36]

324 Here, we found that more than 10% of hallmark-free sequences are in the non-immunized

325 HCAbs repertoire, indicating an increase in the HCAbs diversity by sharing V germline genes

326 with tetrameric IgGs. Interestingly, the germline gene IGHV3S39*01 contributes about 60%
327 of V segment usage among all non-hallmark V germline. The biological mechanism of how
328 hallmark-free HCAbs are developed is still unknown. It is well accepted that Ig heavy chains
329 (HC) are selected at the pre-B cell receptor (pre-BCR) checkpoint. Martin *et al.* found specific
330 structural requirements (CDR3 length and amino composition) to select Ig μ heavy-chains
331 during maturation of the pre-B stage.[39] In our dataset, the hallmark residues (F37, E44, and
332 G47), which usually forming a contact interface with the CDR3 to stabilize the structure,
333 show greater diversity than the others in FR2 except for the ones near CDR regions (Fig. S3
334 in Supplementary Material). Based on our observations, we infer that partial VH germline
335 genes, if not all, are capable of rearranging to HCAbs, but only a small portion (~10%) pass
336 the pre-B checkpoint. Nevertheless, our data confirm that germline gene usage shows a high
337 preference for specific genes. Five out of 17 hallmark-containing germline V segments are
338 responsible for 85.54% of V gene usage in the dataset (Fig.2A).

339   Studies of antibody repertoires from humans and mice demonstrated that germline gene
340 usage is dynamic during vaccination or infection. Hence, we investigated non-immunized
341 samples from three individuals and two samples from one antigen injected animal with two
342 weeks interval. The results show that the V and J germline usage, VJ pairing, DJ pairing, and
343 the amino acid substitution preference are highly correlated whether antigen immunized or
344 not (Table S3 in Supplementary Material). This high similarity is in accordance with a recent
345 work that revealed the high prevalence of shared clonotypes in human B cell repertoires.[40]
346 Contrarily, the D germline usage shows poor correlations especially between the naïve and
347 immunized samples, indicating the D fragment seems to be the main driving force for *in vivo*
348 antibody maturation.

349   The CDR3 is the most polymorphic region of IgGs and the main contributor to antigen
350 binding.[41, 42] The CDR3 loop in VHHs of dromedary is longer than the loop in VHs from
351 humans or mice (average of 14 or 13 AAs, respectively).[43] Longer loop lengths increase the
352 paratope size and consequently help compensate for the diversity loss that occurs when light
353 chains are absent.[9] Analysis of 114 conventional camel antibodies showed that CDR3 length
354 is dependent on the germline gene family.[44] Our HTS data demonstrated that the distribution
355 of CDR3 length various on V germline families, indicating the length of CDR3 loops is
356 determined by the usage of the germline gene. This finding is in accordance with studies of T
357 cell receptor (TCR), whose repertoire distribution patterns depend on the use of the germline
358 genes.[45, 46]

359     Somatic hypermutation (SHM) has long been known as a key process for increasing
360 diversity and improving affinity of antibodies. Comparative analysis of the immune repertoire
361 between the conventional antibodies and the HCAbs from the *Bactrian* camel showed that the
362 nucleotide mutation rate of HCAbs is higher than that of canonical antibodies.[8] In contrast,
363 the calculation of the amino acid mutation rate shows no significant difference in the
364 substitution rate between hallmark and non-hallmark germline gene families (Table S2 in
365 Supplementary Material). However, the substitution patterns of each VJ family did not
366 converge (Fig. 6). Quantitative analysis of GSSPs also confirmed the diversity of the lineages
367 that originated from various VJ germline genes.

368     A recent study on antibody maturation showed that antibodies that respond to the same
369 antigen to a large extent share similar amino acid substitutions.[47] It appears that the conserved
370 antibody structures that drive adaptive immune responses are highly limited and selected.[31]
371 This is consistent with the results of dominant mutations in HCAbs, suggesting the existence
372 of some preferred mutation patterns. For the result of overall ASR, we observed that non-
373 polar amino residues tend to mutate to non-polar AAs; polar and charged residues are more
374 likely to be substituted by polar and charged AAs (Fig.5, a boxed region in heat-map).
375 Phenylalanine (F), alanine (A), serine (S), and aspartic acid (D) are the germline residues that
376 are most preferred to be substituted (Fig.5, upper histogram), while alanine (A), serine (S),
377 glycine (G), aspartic acid (D) are the residues to which are most likely to be mutated (Fig.5,
378 right histogram).

379     The new pipeline developed in this work has revealed novel and detailed features of the
380 HACbs repertoire, which is important for VHH engineering or design. The GSSPs built in this
381 work can describe the mutation sequence space of variable domains of antibodies.[20] In
382 previous studies, we have shown that coupling with appropriate evolutionary profile
383 information, our evolution-based protein design protocol, EvoDesign, exhibits a high
384 accuracy in designing protein folds [48, 49] and protein-protein interactions.[16, 17] The preference
385 of germline usage and mutation of HCAbs can be very useful to reduce the effective searching
386 space of amino acid sequences and increase the accuracy of protein design.

## Acknowledgements

## Authors' contributions

399 ZT and YZ conceived and designed the study. ZT developed the computer code, carried out
400 the statistical analysis and wrote the initial draft. ZT, XH, and YZ reviewed and edited the
401 manuscript. JF, NH, and YL conducted the experiments of sample collection and deep
402 sequencing. WZ participated in the data visualization. All authors read and approved the final
403 manuscript.

## Conflict of Interest

405 The authors have declared no competing interests.

406

## References

408 1. Hamers-Casterman C, Atarhouch T, Muyldermans S, Robinson G, Hamers C, Songa EB,
409     et al. Naturally-occurring antibodies devoid of light-chains. *Nature* 1993; 363:446-8.

410 2. Greenberg AS, Avila D, Hughes M, Hughes A, McKinney EC, Flajnik MF. A new antigen
411     receptor gene family that undergoes rearrangement and extensive somatic diversification in
412     sharks. *Nature* 1995; 374:168-73.

413 3. Herce HD, Schumacher D, Schneider AFL, Ludwig AK, Mann FA, Fillies M, et al. Cell-
414     permeable nanobodies for targeted immunolabelling and antigen manipulation in living
415     cells. *Nat Chem* 2017; 9:762-71.

416 4. Bruce VJ, McNaughton BR. Evaluation of Nanobody Conjugates and Protein Fusions as
417     Bioanalytical Reagents. *Anal Chem* 2017; 89:3819-23.

418 5. Bannas P, Hambach J, Koch-Nolte F. Nanobodies and Nanobody-Based Human Heavy
419     Chain Antibodies As Antitumor Therapeutics. *Front Immunol* 2017; 8:1603.

420 6. Steeland S, Vandenbroucke RE, Libert C. Nanobodies as therapeutics: big opportunities
421     for small antibodies. *Drug Discov Today* 2016; 21:1076-113.

422    7.  Blanc MR, Anouassi A, Ahmed Abed M, Tsikis G, Canepa S, Labas V, et al. A one-step
423         exclusion-binding procedure for the purification of functional heavy-chain and
424         mammalian-type gamma-globulins from camelid sera. *Biotechnol Appl Biochem* 2009;
425         54:207-12.

426    8.  Li X, Duan X, Yang K, Zhang W, Zhang C, Fu L, et al. Comparative Analysis of Immune
427         Repertoires between Bactrian Camel's Conventional and Heavy-Chain Antibodies. *PLoS*
428         *One* 2016; 11:e0161801.

429    9.  Muyldermans S, Baral TN, Retamozzo VC, De Baetselier P, De Genst E, Kinne J, et al.
430         Camelid immunoglobulins and nanobody technology. *Vet Immunol Immunopathol* 2009;
431         128:178-83.

432    10. De Genst E, Silence K, Ghahroudi MA, Decanniere K, Loris R, Kinne J, et al. Strong in
433         vivo maturation compensates for structurally restricted H3 loops in antibody repertoires. *J*
434         *Biol Chem* 2005; 280:14114-21.

435    11. Rubelt F, Busse CE, Bukhari SAC, Burckert JP, Mariotti-Ferrandiz E, Cowell LG, et al.
436         Adaptive Immune Receptor Repertoire Community recommendations for sharing immune-
437         repertoire sequencing data. *Nat Immunol* 2017; 18:1274-8.

438    12. Galson JD, Pollard AJ, Truck J, Kelly DF. Studying the antibody repertoire after
439         vaccination: practical applications. *Trends Immunol* 2014; 35:319-31.

440    13. Ye B, Smerin D, Gao Q, Kang C, Xiong X. High-throughput sequencing of the immune
441         repertoire in oncology: Applications for clinical diagnosis, monitoring, and
442         immunotherapies. *Cancer Lett* 2018; 416:42-56.

443    14. Fridy PC, Li Y, Keegan S, Thompson MK, Nudelman I, Scheid JF, et al. A robust
444         pipeline for rapid production of versatile nanobody repertoires. *Nat Methods* 2014;
445         11:1253-60.

446    15. Turner KB, Naciri J, Liu JL, Anderson GP, Goldman ER, Zabetakis D. Next-Generation
447         Sequencing of a Single Domain Antibody Repertoire Reveals Quality of Phage Display
448         Selected Candidates. *PLoS One* 2016; 11:e0149393.

449    16. Pearce R, Huang X, Setiawan D, Zhang Y. EvoDesign: Designing Protein-Protein Binding
450         Interactions Using Evolutionary Interface Profiles in Conjunction with an Optimized
451         Physical Energy Function. *J Mol Biol* 2019; 431:2467-76.

452    17. Shultis D, Mitra P, Huang X, Johnson J, Khattak NA, Gray F, et al. Changing the
453         Apoptosis Pathway through Evolutionary Protein Design. *J Mol Biol* 2019; 431:825-41.

454  18. Cohen RM, Kleinstein SH, Louzoun Y. Somatic hypermutation targeting is influenced by
455      location within the immunoglobulin V region. *Mol Immunol* 2011; 48:1477-83.

456  19. Cui A, Di Niro R, Vander Heiden JA, Briggs AW, Adams K, Gilbert T, et al. A Model of
457      Somatic Hypermutation Targeting in Mice Based on High-Throughput Ig Sequencing
458      Data. *J Immunol* 2016; 197:3566-74.

459  20. Sheng Z, Schramm CA, Kong R, Program NCS, Mullikin JC, Mascola JR, et al. Gene-
460      Specific Substitution Profiles Describe the Types and Frequencies of Amino Acid Changes
461      during Antibody Somatic Hypermutation. *Front Immunol* 2017; 8:537.

462  21. Tu Z, Xu Y, He QH, Fu JH, Liu X, Tao Y. Isolation and characterisation of
463      deoxynivalenol affinity binders from a phage display library based on single-domain
464      camelid heavy chain antibodies (VHHs). *Food Agric Immunol* 2012; 23:123-31.

465  22. Tu Z, Xu Y, Liu X, He Q, Tao Y. Construction and Biopanning of Camelid NaIve Single-
466      domain Antibody Phage Display Library. *China Biotechnol* 2011; 31:31-6.

467  23. Liu X, Xu Y, Xiong YH, Tu Z, Li YP, He ZY, et al. VHH phage-based competitive real-
468      time immuno-polymerase chain reaction for ultrasensitive detection of ochratoxin A in
469      cereal. *Anal Chem* 2014; 86:7471-7.

470  24. Aronesty E. Comparison of sequencing utility programs. *Open Bioinform Journal* 2013:1-
471      8.

472  25. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: Molecular Evolutionary
473      Genetics Analysis across computing platforms. *Mol Biol Evol* 2018; 35:1547-9.

474  26. Lefranc MP, Giudicelli V, Duroux P, Jabado-Michaloud J, Folch G, Aouinti S, et al.
475      IMGT (R), the international ImMunoGeneTics information system (R) 25 years on.
476      *Nucleic Acids Res* 2015; 43:D413-D22.

477  27. Ye J, Ma N, Madden TL, Ostell JM. IgBLAST: an immunoglobulin variable domain
478      sequence analysis tool. *Nucleic Acids Res* 2013; 41:W34-W40.

479  28. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
480      architecture and applications. *BMC Bioinformatics* 2009; 10:421.

481  29. Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: A sequence logo generator.
482      *Genome Res* 2004; 14:1188-90.

483  30. Arnaout R, Lee W, Cahill P, Honan T, Sparrow T, Weiand M, et al. High-resolution
484      description of antibody heavy-chain repertoires in humans. *PLoS One* 2011; 6:e22365.

485  31. Henry Dunand CJ, Wilson PC. Restricted, canonical, stereotyped and convergent
486      immunoglobulin responses. *Philos Trans R Soc Lond B Biol Sci* 2015; 370.

487    32. Chen L, Kutskova YA, Hong F, Memmott JE, Zhong S, Jenkinson MD, et al. Preferential

488        germline usage and VH/VL pairing observed in human antibodies selected by mRNA

489        display. *Protein Eng Des Sel* 2015; 28:427-35.

490    33. Jayaram N, Bhowmick P, Martin AC. Germline VH/VL pairing in antibodies. *Protein Eng*

491        *Des Sel* 2012; 25:523-9.

492    34. Achour I, Cavelier P, Tichit M, Bouchier C, Lafaye P, Rougeon F. Tetrameric and

493        homodimeric camelid IgGs originate from the same IgH locus. *J Immunol* 2008; 181:2001-

494        9.

495    35. Kabat EA, Wu TT. Identical V region amino acid sequences and segments of sequences in

496        antibodies of different specificities. Relative contributions of VH and VL genes,

497        minigenes, and complementarity-determining regions to binding of antibody-combining

498        sites. *J Immunol* 1991; 147:1709-19.

499    36. Deschacht N, De Groeve K, Vincke C, Raes G, De Baetselier P, Muyldermans S. A novel

500        promiscuous class of camelid single-domain antibody contributes to the antigen-binding

501        repertoire. *J Immunol* 2010; 184:5696-704.

502    37. Monegal A, Olichon A, Bery N, Filleron T, Favre G, de Marco A. Single domain

503        antibodies with VH hallmarks are positively selected during panning of llama (Lama

504        glama) naive libraries. *Dev Comp Immunol* 2012; 36:150-6.

505    38. Flajnik MF, Deschacht N, Muyldermans S. A case of convergence: why did a simple

506        alternative to canonical antibodies arise in sharks and camels? *PLoS Biol* 2011;

507        9:e1001120.

508    39. Martin DA, Bradl H, Collins TJ, Roth E, Jack HM, Wu GE. Selection of Ig mu heavy

509        chains by complementarity-determining region 3 length and amino acid composition. *J*

510        *Immunol* 2003; 171:4663-71.

511    40. Soto C, Bombardi RG, Branchizio A, Kose N, Matta P, Sevy AM, et al. High frequency

512        of shared clonotypes in human B cell receptor repertoires. *Nature* 2019; 566:398-402.

513    41. De Genst E, Saerens D, Muyldermans S, Conrath K. Antibody repertoire development in

514        camelids. *Dev Comp Immunol* 2006; 30:187-98.

515    42. Miqueu P, Guillet M, Degauque N, Dore JC, Soulillou JP, Brouard S. Statistical analysis

516        of CDR3 length distributions for the assessment of T and B cell repertoire biases. *Mol*

517        *Immunol* 2007; 44:1057-64.

518    43. Kaplinsky J, Li A, Sun A, Coffre M, Koralov SB, Arnaout R. Antibody repertoire deep
519        sequencing reveals antigen-independent selection in maturing B cells. *Proc Natl Acad Sci*
520        *U S A* 2014; 111:E2622-9.

521    44. Griffin LM, Snowden JR, Lawson AD, Wernery U, Kinne J, Baker TS. Analysis of heavy
522        and light chain sequences of conventional camelid antibodies from Camelus dromedarius
523        and Camelus bactrianus species. *J Immunol Methods* 2014; 405:35-46.

524    45. Nishio J, Suzuki M, Nanki T, Miyasaka N, Kohsaka H. Development of TCRB CDR3
525        length repertoire of human T lymphocytes. *Int Immunol* 2004; 16:423-31.

526    46. Gomez-Tourino I, Kamra Y, Baptista R, Lorenc A, Peakman M. T cell receptor beta-
527        chains display abnormal shortening and repertoire sharing in type 1 diabetes. *Nat Commun*
528        2017; 8:1792.

529    47. Tian M, Cheng C, Chen X, Duan H, Cheng HL, Dao M, et al. Induction of HIV
530        Neutralizing Antibody Lineages in Mice with Diverse Precursor Repertoires. *Cell* 2016;
531        166:1471-84.

532    48. Brender JR, Shultis D, Khattak NA, Zhang Y. An Evolution-Based Approach to De Novo
533        Protein Design. *Methods in molecular biology (Clifton, N.J.)* 2017; 1529:243-64.

534    49. Mitra P, Shultis D, Zhang Y. EvoDesign: De novo protein design based on structural and
535        evolutionary profiles. *Nucleic Acids Res* 2013; 41:W273-80.

536    50. John Towns, Timothy Cockerill, Maytal Dahan, Ian Foster, Kelly Gaither, Andrew
537        Grimshaw, et al. XSEDE: Accelerating Scientific Discovery. *Comput Sci Eng* 2014; 16:62-
538        74.

539 **Table 1. Summary of sequencing datasets in this work.**

| Sample [a] | Data [b] (counts) | Joined Data [c] (counts) | Unique DNA Data (counts) | Unique AA Data (counts) |
|---|---|---|---|---|
| Naïve-1 | 14,130,770 × 2 | 5,850,191 | 2,550,586 | 1,973,186 |
| Naïve-2 | 9,783,739 × 2 | 6,381,831 | 2,317,312 | 1,717,133 |
| Naïve-3 | 7,939,987 × 2 | 5,285,912 | 1,763,675 | 1,271,695 |
| Immune-1 | 2,935,298 × 2 | 1,841,952 | 1,270,186 | 782,529 |
| Immune-2 | 3,459,366 × 2 | 2,270,196 | 1,653,673 | 1,149,386 |

540 a. Naïve-1, Naïve-2, and Naïve-3 were collected from three healthy donors; Immune-1
541 and Immune-2 were collected from one donor after fifth and seventh immunization,
542 respectively.
543 b. Data is the total sequences of paired-end reads after filter and quality control;
544 c. Joined Data is the number of sequences that were generated from paired-end reads.
545

546

547 **Figure Legends**

548

549 **Figure 1. Workflow for the analysis of HCAbs repertoire.** The coding sequences of VHHs

550 are amplified from long and short hinge HCAbs, respectively, where the next-generation

551 sequencing (NGS) data are processed using fastq-join to merge paired-end sequences after

552 discarding the low-quality reads. Next, IgBLAST is used to assign V(D)J genes for each

553 transcript. Coding sequence (CDS) length distribution, V(D)J usage, VJ pairing, and DJ

554 pairing usage are based on DNA sequences, while other analyses are based on amino acid

555 sequences.

556

557 **Figure 2. V, D, and J germline gene usage of VHHs repertoire is highly biased.** (A) Usage

558 of V germline genes. (B) Usage of D germline genes. (C) Usage of J germline genes.

559

560 **Figure 3. V(D)J pairing of VHHs is biased.** (A) VJ pairing of germline genes in naïve

561 VHHs repertoire is highly biased to pairs of certain germline genes. The top 21 VJ pair made

562 up > 90% of the clones in the repertoire. (B) and (C) are the *in silico* simulation of DJ and VJ

563 combination. The error bars illustrate the relative deviation of the 2,000 steps of simulation,

564 while the columns represent the relative deviation of the VHHs repertoire.

565

566 **Figure 4. Distribution of CDR3 amino acid length is depending on germline gene**

567 **families.** The average CDR3 lengths of the top 11 V germline genes are presented in the

568 upper right panel. Hallmark and non-hallmark residue V genes are shown in red or blue bars,

569 respectively.

570

571 **Figure 5. Most substitutions are rare and partial mutations are preferred.** Due to

572 mutation bias, most observed types of substitution occur rarely. The histogram in the upper

573 and right are the sum of the cases in the corresponding column or row, respectively. The

574 asterisk in mutation type means stop codon. Dashes mean gaps in both germline residues and

575 mutation types.

576

577 **Figure 6. The similarity of GSSPs between VJ germline gene families.** The Jensen-

578 Shannon divergence was used to compare GSSPs, and the distance matrix was visualized

579    using multidimensional scaling. The VJ pair types that have more than 100 (A and B), 500 (C

580    and D), or 1000 (E and F) lineages were calculated and plotted respectively. GSSPs of the

581    same V gene tend to be clustered together. Meanwhile, the distance of partial VJ pair is close,

582    which indicates sharing similar mutation patterns. VJ pairs from the top 11 V genes are shown

583    in the right column (B, D, and F).

584

585    **Figure 7. Comparison of V(D)J usage of long-hinge (IgG2) and short-hinge (IgG3)**

586    **HCAbs.** Sequences belonging to IgG2 or IgG3 were identified using specific primers for

587    BLAST. (A), (B), and (C) are the V, D, and J gene usage of both types of HCAbs,

588    respectively.

imm_13224_f1.jpg

imm_13224_f2.jpg

imm_13224_f3.jpg

imm_13224_f4.jpg

imm_13224_f5.jpg

imm_13224_f6.jpg

imm_13224_f7.jpg