# Large-scale Analysis of 2,152 dataset reveals key features of B cell biology and the antibody repertoire

Xiujia Yang[1,2,a*], Minhui Wang[1,b*], Dianchun Shi[3*], Yanfang Zhang[2,c*], Huikun Zeng[1,2*], Yan Zhu[2*],

Chunhong Lan[1,2,d*], Jiaqi Wu[2], Yang Deng[4], Shixin Guo[5], Lijun Xu[6], Cuiyu Ma[1,2], Yanxia Zhang[2],

Rongrong Wu[3], Jinxia Ou[7,e], Chu-jun Liu[8], Changqing Chang[9], Wei Yang[10,f], Huijie Zhang[11], Jun

Chen[12], Lijie Qin[6], Hongwei Zhou[7,g], Jin-Xin Bei[8], Lai Wei[5], Guangwen Cao[4†], Xueqing Yu[3†],

Zhenhai Zhang[1,2,13,14,h†]

9    [1]State Key Laboratory of Organ Failure Research, National Clinical Research Center for Kidney

10   Disease, Division of Nephrology, Nanfang Hospital, Southern Medical University, Guangzhou,

11   510515, China

12   [2]Department of Bioinformatics, School of Basic Medical Sciences, Southern Medical University,

13   Guangzhou 510515, China

14   [3]Department of Geriatrics, Guangzhou First People's Hospital, School of Medicine, South China

15   University of Technology, Guangzhou 510030, China

16   [4] Department of Epidemiology, Second Military Medical University, 800 Xiangyin Rd., Shanghai

17   200433, China

18   [5]State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-Sen University,

19   Guangzhou 510060, China

20   [6]Department of Emergency, Henan Provincial People's Hospital, Zhengzhou 450003, China

21   [7]Microbiome Medicine Center, Division of Laboratory Medicine, Zhujiang Hospital, Southern

22   Medical University, Guangzhou 510282, China

23   [8]Sun Yat-Sen University Cancer Center, State Key Laboratory of Oncology in South China,

24   Collaborative Innovation Center for Cancer Medicine, Sun Yat-Sen University, Guangzhou 510030,

25   China

26   [9]Integrate Microbiology Research Center, South China Agricultural University, Guangzhou, 510642,

27   China

28   [10]Department of Pathology, School of Basic Medical Sciences, Southern Medical University,

29   Guangzhou, 510515, China

30   [11]Department of Endocrinology and Metabolism, Nanfang Hospital, Southern Medical University,

31   Guangzhou 510515, China

32    [12]MOE Laboratory of Biosystems Homeostasis & Protection and Innovation Center for Cell Signaling

33    Network, College of Life Sciences, Zhejiang University, Hangzhou 310058, China

34    [13]Center for precision medicine, Guangzhou First People's Hospital, School of Medicine, South China

35    University of Technology, Guangzhou 510030, China

36    [14]Key Laboratory of Mental Health of the Ministry of Education, Guangdong-Hong Kong-Macao

37    Greater Bay Area Center for Brain Science and Brain-Inspired Intelligence, Southern Medical

38    University, Guangzhou 510515, China

39    [*]These authors contributed equally to this work.

40    [†]To whom correspondence should be addressed:

41    Zhenhai Zhang, zhenhaismu@163.com

42    Xueqing Yu, yuxueqing@gdph.org.cn

43    Guangwen Cao, gcao@smmu.edu.cn

44

45

46    ORCID:

47    [a]0000-0003-4036-4995                    [b]0000-0001-8121-7786

48    [c]0000-0001-9309-7347

49    [d]0000-0001-5030-8247                    [e]0000-0002-1680-2425

50    [f]0000-0001-9438-7215

51    [g]0000-0003-2472-8541                    [h]0000-0002-4310-0525

52

53    *Running title: Key features of antibody repertoire.*

# 54    **Abstract**

55    Antibody repertoire sequencing (Ig-seq) has been widely used in studying humoral responses, with

56    promising results. However, the promise of Ig-seq has not yet been fully realized, and key features of

57    the antibody repertoire remain elusive or controversial. To clarify these key features, we analyzed

58    2,152 high-quality heavy chain antibody repertoires, representing 582 donors and a total of 360

59    million clones. Our study revealed that individuals exhibit very similar gene usage patterns for

60    germline V, D, and J genes and that 53 core V genes contribute to more than 99% of the heavy chain

61    repertoire. We further found that genetic background is sufficient but not necessary to determine usage

62    of V, D, and J genes. Although gene usage pattern is not affected by age, we observed a significant

63    sex preference for 24 V genes, 9 D genes and 5 J genes, but found no positional bias for V-D and D-J

64    recombination. In addition, we found that the number of observed clones that were shared between

65    any two repertoires followed a linear model and noted that the mutability of hot/cold spots and single

66    nucleotides within antibody genes suggested a strand-specific somatic hypermutation mechanism.

67    This population-level analysis resolves some critical characteristics of the antibody repertoire and thus

68    may serve as a reference for research aiming to unravel B cell-related biology or diseases. The metrics

69    revealed here will be of significant value to the large cadre of scientists who study the antibody

70    repertoire.

71    **Keywords**: B-cell biology, antibody repertoire, large-scale analysis, high-throughput sequencing,

72    Ig-seq

73

74

75

76

# Introduction

The antibody repertoire is defined as the entire collection of B-cell receptors and antibodies that grant protection against a plethora of pathogens. A deeper understanding of the antibody repertoire under normal physiological conditions and in pathogenic conditions may shed light on functional immune responses and reveal the full scope of their protective and pathogenic functions. However, despite this great potential, collecting enough antibody molecules to capture the immense diversity of the antibody repertoire has been a critical challenge.

Using high-throughput sequencing technology, Weinstein et al. developed antibody repertoire sequencing (Ig-seq) (Weinstein et al., 2009), which allows researchers to capture millions or even billions of antibody variable regions within a single experiment. The vast amount of data acquired by Ig-seq enables a deeper and more thorough evaluation of the key features of the antibody repertoire, as well as its constituent antibody molecules, at the single-nucleotide level. In the past decade, Ig-seq has advanced the study of many important sub-fields of B-cell immunology, such as antibody discovery (Reddy et al., 2010; Zhu et al., 2013a), vaccination development (Jackson et al., 2014; Jiang et al., 2013; Joyce et al., 2016; Li et al., 2012), infection (Krebs et al., 2019; Parameswaran et al., 2013; Wu et al., 2015; Wu et al., 2011a), allergy (Hoh et al., 2016; Patil et al., 2015; Wu et al., 2014), autoimmune disease (Stern et al., 2014; Tipton et al., 2015; von Büdingen et al., 2012), and cancer immunology (Faham et al., 2012; Gawad et al., 2012; Kurtz et al., 2015). For example, using Ig-seq coupled with single-cell cloning technology, we and others identified thousands of HIV-1-neutralizing antibodies that bind to different epitopes and delineated their lineage-dependent maturation pathways (Bonsignori et al., 2016; Wu et al., 2015; Wu et al., 2011b; Zhu et al., 2013b). Studies of antibody repertoires after virus infection also led to the discovery of antibody convergence – a mechanism whereby identical or very similar antibody clonotypes are generated in different individuals facing the same selective pressure (Parameswaran et al., 2013). These results suggested that the antibody

101    repertoire could be used to track an individual's immune history as well as to monitor the

102    immunological memory of a community.

103        The use of antibody repertoire in autoimmune diseases has provided important insight into both

104    disease mechanisms and fundamental B cell biology. For example, Tipton et al. revealed that systemic

105    lupus erythematosus (SLE) autoreactivity occurred during a polyclonal activation of

106    IGHV4-34-dominant B cell clones via both germinal center-dependent and germinal

107    center-independent mechanisms (Tipton et al., 2015). Büdingen et al. discovered that a pool of clonal

108    related antibodies participates in a robust bidirectional exchange across the blood-brain barrier (von

109    Büdingen et al., 2012). Analyzing the antibody repertoires of patients with the same disease, Stern et

110    al. found that majority of the disease-related autoantibodies matured outside of the central nervous

111    system and trafficked freely across tissue barriers (Stern et al., 2014).

112        Prior studies using Ig-Seq accumulated a wealth of antibody repertoire data. This under-explored

113    population-level big data could potentially help us resolve the important yet unclear or controversial

114    features of the B cell biology and the antibody repertoire. For instance, what are the germline V, D,

115    and J gene usage patterns and how similar are they between individuals? What are the factors

116    determining these patterns if they do exist? Is there preferential recombination between V-D and D-J

117    genes? What are the rules that govern the somatic hypermutations (SHM)? What are the proportion of

118    public clones between individuals and what functions do these clones exert?

119        With these unsolved or controversial questions in mind, we collected 2,152 high quality antibody

120    heavy chain repertoires and performed thorough and in-depth analyses. These analyses revealed

121    patterns of B cell biology as well as key features of the antibody repertoire, which will be of

122    significant value to the large cadre of scientists in the field.

123

# Results

## Overview of datasets used in this study

The immense diversity of B cells is derived from two important biological processes: germline gene segment recombination, which introduces indels in complementarity-determining region 3, and activation-induced cytidine deaminase, which leads to somatic hypermutations in the antibody variable regions during affinity maturation. Thus, any mutation or indel in the variable region may be important for understanding and revealing B cell biology. It is thus essential to understand the sequencing errors that are intrinsic to Ig-seq approaches. For example, 454 sequencing often generates indels in the homopolymer region, and PCR amplification and high-throughput sequencing can also generate base errors and chimeras. These intrinsic errors would be easily mistaken as somatic hypermutation (SHM) generated during affinity maturations. We therefore only included samples that were sequenced by Illumina instruments. We also required the sequencing reads to cover a minimum of 500 bp, the primers to capture the full spectrum of antibodies generated by any V(D)J recombination, and a minimum number of productive reads (Materials and Methods). After filtering on these stringent criteria, we identified a total of 1,857 repertoires from 33 published studies and 295 repertoires from in-house sequencing for further analysis (Figure 1a).

The sample-associated metadata, including age, sex, physiological condition, tissue origin, and amplification method, are shown in Figures 1b-f. Age composition was more balanced for the sampled individuals than for the samples (Figure 1b and Figure S1), and the number of individuals for each age group is more than 30 (more than 80 samples for each age group). Slightly more than half of the samples were from females (Figure 1c). Sequencing libraries for all recruited samples were mainly amplified using multiplex PCR (Figure 1d). Donor conditions and the sources of the samples were classified into 13 and 6 directories respectively (Figures 1e and 1f). These repertoires covered a broad spectrum of diseases, such as autoimmune disease, cancer, virus infection, and more. The majority

148    (76.8%) of samples were derived from peripheral blood mononuclear cells (PBMCs); there were also

149    samples from bone marrow, intestine, lung, and spleen. Overall, we included a total of 7,378,354,271

150    raw reads in our analysis.

## The core V gene set determines the clear majority of antibodies

152    The variable usage of germline genes represents the first level of antibody repertoire diversity

153    and is believed to affect immune function (Glanville et al., 2011). Naïve germline gene usage may be

154    optimized for interactions with common antigens and may serve as a control to detect

155    pathology-driven repertoire variation in the B cell memory compartment (Laserson et al., 2014). For

156    these reasons, the gene usage pattern has been studied at a small scale and under different

157    experimental settings using two different quantification methods: gene usage and gene expression.

158    Gene usage quantifies genes at the level of individual clones, whereas gene expression quantifies the

159    occurrence of genes with each read. Gene expression is sensitive to cell type composition, such as

160    clonal expansion in response to an adaptive immune stimulus, and thus is less optimal for comparisons

161    between samples with differences in source tissues, immune status, and donor health.

162    Library preparation technique also affects the quantification of genes. Two amplification

163    strategies were used in the high-quality Ig-seq datasets: multiplex polymerase chain reaction (MPCR)

164    and rapid amplification of cDNA ends (RACE). Previous studies showed that MPCR can introduce

165    bias in library sequencing, even with an optimized primer set, while RACE introduces less bias (He et

166    al., 2015; Liu et al., 2016; Robins, 2013). We compared quantitative metrics for 1,409 and 743

167    samples amplified by MPCR and RACE, respectively. D and J gene usage were less influenced by

168    either RACE or MPCR. However, V gene usage was more consistent between RACE and MPCR

169    (Figure S2a-f and Materials and Methods) due to the various primers on the 5' ends. We therefore

170    selected gene usage for the following analyses unless otherwise specified.

171    It has been long known that V(D)J gene segments are preferentially used or expressed, and the

172    idea of a core gene set has been proposed (Boyd et al., 2010). However, which genes are in the core

173    set and the extent to which they contribute to the antibody repertoire remains unclear. Taking

174    advantage of the large data set used in this study, we plotted the gene usage of V, D, and J genes. As

175    shown in Figure 2, although the number of V genes present in each sample varies from 15

176    (SRR3620039) to 99 (SRR8365259 and SRR4417619) along with the sequencing depth, we observed

177    preferential usage of V genes. For example, IGHV3-30 and IGHV3-23 were present in all samples,

178    while IGHV3-30-52 only appeared in one sample (SRR4417619). Accounting for both the prevalence

179    of specific genes and their contribution to the antibody repertoire, we identified a core set of 53 V

180    genes (Figures 2b and 2c, Materials and Methods). To our surprise, there are 3 pseudogenes,

181    IGHV3-11 (2,147 samples, 581 donors), IGHV3-69-1 (2,128 samples, 579 donors), and IGHV3-71

182    (1,691 samples, 464 donors), in the core gene set. All core V genes contribute to a median of 99.33%

183    of clones (Figure S3). The remaining V genes thus either contribute little to the repertoire or are not

184    present. IGHJ3, IGHJ4, and IGHJ5 are present in all samples, while IGHJ1, IGHJ2, and IGHJ6 occur

185    in 2,149, 2,150, and 2,151 samples, respectively. Of these, IGHJ4 and IGHJ6 are found in a median of

186    50.56% and 18.37% clones. IGHJ1 is the least used gene, contributing to a median of 2.15% clones.

187    Three D genes, IGHD3-10, IGHD3-22, and IGHD6-19, are more prevalent. Statistics for V, D, and J

188    gene segments are shown in Table S1, sorted by their occurrence in 582 individuals.

189    **Genetic background is sufficient but not necessary for achieving consistent germline**

190    **gene usage patterns**

191        The factors that determine V gene usage patterns have been of great interest in the field, with

192    different studies yielding different results. By comparing the repertoires of monozygotic twins and

193    unrelated individuals, Glanville et al. concluded that gene usage patterns are heritable, whereas

194    Arnaout et al., Briney et al., and Laserson et al. reported that an individual's gene usage pattern is

195    almost identical or remarkably consistent among individuals (Briney et al., 2012; Glanville et al., 2011;

196    Laserson et al., 2014). Thus, the effect of genetic background on V gene usage pattern is still unclear.

197        We therefore selected 109 repertoires, all amplified using 5'RACE, from 23 unrelated males and

198    3 pairs of monozygotic twins. From these repertoires, we calculated Pearson's correlation coefficients

199  (Pearson's r values) for pairwise V gene usage (sample pairs from the same donor were excluded). As

200  shown in Figure 3a, the overall coefficients of all 104 genes are higher than 53 core genes. Further

201  scrutinizing the data revealed that most of the non-core genes had values of 0 (Figure S4a and S4b).

202  For the male-derived samples, the minimum and maximum number of uncaptured core genes are 0

203  and 13, respectively, with a median of 1 and a mean of 1.53. For the non-core gene set, the minimum

204  number is 5, with a maximum of 47, a median of 21, and a mean of 23.12. For the female-derived

205  samples, the minimum and maximum number of uncaptured core genes were 0 and 12, respectively,

206  with a median of 1 and a mean of 1.10. For the non-core gene set, the minimum number is 5, with a

207  maximum of 43, a median of 21, and a mean of 22.10. These values elevated the pairwise coefficients.

208  Therefore, we decided to use the 53 core V genes identified earlier (Figures 2b and 2c) for further

209  analyses. The Pearson's r values of unrelated donors ranged from 0.3681 to 0.9517, while those of

210  monozygotic twins of the same cell type ranged from 0.9130 to 0.9952. The higher coefficient

211  observed in monozygotic twins indicates that genetic background is sufficient to account for

212  consistent V gene usage. However, we also observed 16 unrelated sample pairs that showed a

213  coefficient higher than 0.9130, the minimum coefficient observed between monozygotic twins with

214  the same cell type. Thus, a shared genetic background is not necessary for generating repertoires with

215  very similar V gene usage. The usage patterns of D and J genes also showed the same phenomena

216  (Figure S5a and S5b), and results were similar in the female-derived samples (Figure S5c-e). We

217  therefore conclude that genetic background plays a critical role in defining antibody repertoire by

218  influencing germline gene usage. However, individuals with different genetic backgrounds can also

219  achieve a remarkably similar repertoire.

220  **V, D, and J gene usage shows sex and isotype preferences**

221  After defining the relationship between gene usage and genetic background, we went on to

222  analyze two other major factors: age and sex. Consistent with a previous study (Wang et al., 2014),

223  our results showed that there is no linear relationship between gene usage and age, regardless of sex

224  (Figures 3b and Figure S6a-f). To rigorously examine the impact of sex differences on antibody

225  repertoire gene usage, we calculated pairwise germline gene usage patterns for 499 healthy PBMC

226  samples amplified with RACE from 94 male and 164 female individuals. To our surprise, we observed

227  that 24 core V genes (Figure 3c), 5 J genes (Figure 3d), and 9 D genes (Figure 3e) showed significant

228  differences between male- and female-derived samples (p < 0.05, P values are listed in Table S2).

229      Previous studies reported that patients with influenza and SLE had characteristic changes in

230  antibody gene usage (Pugh-Bernard et al., 2001; Sui et al., 2009). We therefore used healthy donors as

231  background and investigated gene usage in individuals with different diseases (Figure S7a-f and Table

232  S3). For the female-derived samples, 12 V and 8 D genes were out of the range defined by 283 healthy

233  samples. We found that IGHV4-38-2 (SRR4026039 and SRR4026040) and IGHV3-23D

234  (SRR4026032, SRR4026025, and SRR4026031) had increased usage in 1 and 2 out of 6 female

235  Myasthenia Gravis patients. IGHD4-17 (SRR4026038 and SRR7230358), and IGHD3-3

236  (SRR4026022 and SRR4026031) were upregulated in 2 out of 6 female Myasthenia Gravis patients.

237  For the male-derived samples, 19 V, 7 D, and 3 J genes had either higher or lower usage compared to

238  216 healthy male samples. For example, IGHV1-18 (H7N9_00004 and H7N9_00009) and IGHV3-73

239  (H7N9_00011 and H7N9_00005) showed higher usage in 2 out of 4 H7N9-infected samples. Thus, a

240  statistical analysis of large data sets may be a powerful tool in studying the antibody repertoires of

241  unhealthy individuals.

242      We also examined gene usage in different antibody isotypes, namely IgA, IgD, IgG, and IgM.

243  There were total 51 repertoires from 5 females (14 samples) and 12 males (37 samples) available for

244  this analysis (Figure S8). IgA and IgG were clustered together, while IgD and IgM gathered in the

245  same subtree within the same donor. This is true for male (Figure S9a-c) and female (Figure S10a-c)

246  samples and is consistent with a previous study (Laserson et al., 2014).

247  **DJ recombination shows no positional bias**

248      During the recombination process, exonuclease trimming and the random addition of nucleotides

249  between VD and DJ segments create diverse junctions to account for a substantial amount of antigens

250  that may be encountered (Early et al., 1980; Tonegawa, 1983). These junctions, together with the D

251     genes, are known as complementarity-determining region 3 (CDR3), which largely determines the

252     binding specificity of an antibody (Chothia et al., 1989). Due to the functional importance of CDR3,

253     there have been extensive studies looking at VDJ recombination preferences and indels in the

254     junctions (Hansen et al., 2015; Hong et al., 2018; Saada et al., 2007; Souto-Carneiro et al., 2005;

255     Truck et al., 2015).

256     For the recombination bias studies, D and J gene segments were first classified as 5D, 3D, 5J, or

257     3J based on their position on the chromosome. The 5D and 5J categories include the D and J segments

258     located in the upstream region of their respective cluster. The 3D and 3J categories include the

259     downstream D and J gene segments. Thus, 3D and 5J segments are proximal, and 5D and 3J segments

260     are distal (Hong et al., 2018; Saada et al., 2007; Souto-Carneiro et al., 2005; Truck et al., 2015).

261     Comparing DJ recombination between neonates and adults, Souto et al. found that 3D segments

262     preferentially coupled to 5J segments (a proximal bias) throughout development, while 5D segments

263     showed biased recombination to 3J segments (a distal bias) in full-term neonates rearrangements

264     (Souto-Carneiro et al., 2005). Kidd et al. also observed a clear recombinational preference of 5D to 3J

265     and 3D to 5J segments (Kidd et al., 2016). We thus plotted VD and DJ recombination (Figure 4a)

266     using a total of 352 million productive clones that have D genes assigned (Figure 4b). Surprisingly,

267     apart from the preferential usage of core V genes, IGHJ4, IGHJ6, and a few IGHD genes, we did not

268     observe either proximal DJ or distal DJ recombination biases in our data. However, the datasets from

269     previous neonate donors did not meet our inclusion criteria, so we cannot evaluate the positional bias

270     of DJ recombination during neonatal development.

271     **D-D fusion exhibits isotype and distance preferences**

272     D-D fusion, the incorporation of multiple diversity (D) genes during heavy chain recombination,

273     contributes markedly to antibody repertoire diversity and has been thought to generate long CDR3

274     loops that frequently associate in self-reactive and polyreactive antibodies (Briney et al., 2012;

275     Larimore et al., 2012). Briney et al. reported the first quantification of V(DD)J recombinants in naïve,

276     memory IgM and IgG B cells from peripheral blood using Roche 454 sequencing of 4 healthy donors

277    (Briney et al., 2012). Using stringent criteria, they found no antibodies with D-D fusion in the memory

278    IgG population. They also reported that D gene order in cases of D-D fusion matches the order of their

279    loci in the genome.

280        In bulky antibody repertoire sequencing, it is common practice to use different 3' primers

281    targeting different isotypes. We therefore went on to explore D-D fusion in different repertoires as

282    well as different isotypes using IgScout (Safonova and Pevzner, 2019). We first examined CDR3

283    length. Total CDR3s displayed a normal distribution with a peak length of 48 nucleotides. However,

284    the lengths of CDR3s with D-D fusions were much longer, with a peak length of 66 nucleotides

285    (Figure 4b). Hence, D-D fusion does result in longer CDR3s.

286        To explore how often D-D fusion recombinants present in different antibody isotypes, we chose

287    repertoires with at least 5,000 C gene assigned clones for corresponding isotype and calculated the

288    frequency of D-D fusions. IgD had the highest D-D fusion frequency of 0.260% (median, n=104),

289    followed by 0.216% of IgM (median, n=594). IgG (median, n=489) and IgA (n=163) exhibited much

290    lower D-D fusion frequencies of 0.089% and 0.060% (median value), respectively (Figures 4c and

291    Table S4). We did not calculate the D-D fusion frequencies for IgE because too few repertoires were

292    available. These results are consistent with previous findings that D-D fusion recombinants may be

293    negatively selected during isotype switching (Souto-Carneiro et al., 2005).

294        In contrast with previous findings, however, gene order in D-D fusion did not match the order of

295    the corresponding loci in the genome (Figure 4d). The upstream D gene is defined as the "first" D

296    gene (D1) in the fused recombinants and the downstream D genes could be the "second" D gene (D2).

297    In other word, the first D gene (D1) located more 5' in the genome prefer to be the second D gene (D2)

298    in a D-D fusion event. However, D1 gene seem to prefer to fuse with downstream D genes with a span

299    of 7 (The span of adjacent D gene is 1) (Figure 4e). Surprisingly, we did not observe a positive

300    correlation between D-D fusion and D gene usage (Briney et al., 2012). The most abundant pairs were

301    D3-10-D1-1 and D6-1-D1-1. D6-19 often served as D2, and D5-12 or D6-13 as D1. These findings

302    may shed light on the recombination mechanistic studies.

## Stochastic recombination contributes to the public clone

Public clones are defined as to clonotypes shared by multiple individuals (Greiff et al., 2015; Jackson et al., 2013; Miho et al., 2019). It has been suggested that public clones are valuable for designing vaccines, monitoring the immune response to infection or vaccination, developing biomarker patterns of disease states, and mediating the undesirable immune responses associated with autoimmune diseases (Briney et al., 2019; Bürckert et al., 2017; Greiff et al., 2017; Maecker et al., 2012). Recent studies reported that individuals exposed to the same antigen, such as HIV, influenza, or dengue, may develop identical or similar Ig sequences – a phenomenon called antibody convergence (Jackson et al., 2014; Parameswaran et al., 2013; Setliff et al., 2018; Truck et al., 2015). Thus, a comprehensive atlas of public clones may help reconstruct the immunological history of an individual and may enable immunotherapeutic targeting within a population with a specific disease.

Ig-seq has enabled public clone studies via multiple means. Greiff et al. developed an approach that learned the high-dimensional immunogenomic features from the repertoire and enabled the prediction of public and private clones (Greiff et al., 2017). By comparing multiple donors' ultra-deep repertoire sequencing data, Burton et al. and Soto et al. estimated the fraction of public clones in an individual to be approximately 1% and 1% to 6%, respectively (Briney et al., 2019; Soto et al., 2019). Taking advantage of the unprecedented amount of data collected for the present study, we investigated the prevalence of public clones in 2,152 samples. As shown in Figure 5a, we found that the abundance of public clones in a sample decreased when the total number of clones decreased. This result suggests that methodological undersampling may compromise the detection of public clones (Greiff et al., 2015). Furthermore, we also found that the number of public clones in two samples correlates linearly with the product of their respective clone numbers (Figure 5b) and that this correlation improves when the clone numbers for both samples increase (Figure S11a-g). The total number of clones in a given volume of blood varies with an upper boundary. Thus, getting more clones requires more blood samples. Based on different methods, the total clones in an individual's circulating blood has been estimated to be between twenty-five million and one billion (Briney et al., 2019; Soto et al., 2019).

329    Using our linear models with a minimum number of 5 million clones, an individual may possess

330    between $4.2 \times 10^4$ and $6.87 \times 10^7$ public clones in his or her circulating blood.

331         More in-depth analyses revealed that V and J gene usage is almost identical to the gene usage for

332    all clones (Figure S12). On the other hand, public clones possess significantly shorter CDR3s (Figure

333    S13). Statistical analyses of the deletions, non-template additions, and P additions showed significant

334    differences in most elements between public clones and private clones (Figure S14). In particular, the

335    N1 and N2 additions in public clones between the VD and DJ junctions were shorter than those of

336    private clones. This may explain the short CDR3s in public clones and why D genes could not be

337    assigned in many public clones (Figure S15).

338         Of the 162,975 (transformed to the number of unique CDR3 amino acid sequences) public clones

339    identified in this study, 1,059 CDR3s were identical to published antigen-specific or

340    disease-associated antibodies (Figure S16a). Further analyses showed that these CDR3s are enriched

341    for the HIV, influenza, hematological malignancies, EBV, tetanus, and rheumatic categories (Table

342    S5). This enrichment confirmed that antibody convergence was a source of public clones. In addition

343    to the CDR3 enrichment in the antibodies with rheumatic autoimmune disease, we also found a CDR3

344    corresponding to SLE-specific antibodies in one of the healthy donors in our data. In addition, the

345    clonotypes shared by more donors were more abundant (Figure 5c), and this change in abundance was

346    not related to CDR3 length. Previous studies in T cell receptors (TCRs) found that shared TCRs are

347    more likely to be autoreactive (Madi et al., 2014) and that these autoreactive TCRs are important for

348    maintaining an individual's health. It is possible that public clones in an antibody repertoire serve the

349    same function. Surprisingly, we also found 31,226 (66.2%) IgM and 7,699 (70.7%) IgD clones with

350    identical sequences compared to their respective germline V and J genes (see Materials and Methods).

351    These clones are generated solely by VDJ recombination but have no somatic hypermutation,

352    regardless whether the individuals have been exposed to antigen or not. This result suggests that in

353    addition to antibody convergence, the stochastic nature of somatic recombination alone could be a key

354 mechanism of generating public clones. We believe this collection of public clones will be helpful for

355 studies relating to vaccine and therapeutic design targeting shared antibodies.

## Strand specificity features somatic hypermutation

357 Somatic hypermutation (SHM) takes place in the germinal centers of peripheral lymphoid tissues

358 and increases the number of realizable antibodies by several orders of magnitude in addition to

359 combinatorial diversity. The preferences and patterns of SHM allow us to trace the clonal evolution of

360 antibodies under the selective pressure of particular antigen and to facilitate vaccine design (Schramm

361 and Douek, 2018). The nucleotides and amino acid sites that are preferred or disfavored in SHM have

362 been investigated using limited data and in model systems (Schramm and Douek, 2018). SHM in the

363 antibody repertoire results from two types of sequential events. First, activation-induced cytidine

364 deaminase and other molecular components of the SHM machinery introduce mutations to the

365 antibody variable regions. The selective pressure of a particular antigen then acts on these mutations

366 and preserves the favored ones. Thus, the majority of SHM studies worked from unproductive reads to

367 emphasize the effect of mutations rather than the effect of antigen selection. This is particularly

368 beneficial for mechanistic research on SHM because it simplifies the model. However, antigens only

369 place selective pressure on antibodies containing mutations and do not introduce additional mutations.

370 Despite some antigens that may preferably retain rare mutations, the clear majority of mutations in

371 functional antibodies would also reflect the selective flavors of the SHM machinery. Moreover, the

372 selective bias only acts on antigen-specific clones. Thus, this bias would be minimized if the effects of

373 clonal expansion are compensated or removed during computational analyses.

374 With this in mind, we used consensus sequences and a position weight matrix (PWM) to

375 represent a clone and probed mutations in the V genes at different levels (Materials and methods). We

376 first depicted the mutational propensities at the single nucleotide level (Figures S17a and S17b). Three

377 types of transitions (*A* to *G*, *G* to *A*, and *C* to *T*) occurred with high frequencies, and *T* to *C* mutations

378 occurred at a significantly lower frequency. Transversions between purines and pyrimidines were less

379 frequent, with the exception of *G* to *C* mutations. While *C* and *G* showed comparable mutability, *A*

380  exhibited significantly higher mutability than T. Because *A*: *T* and *G*: *C* present as pairs on the

381  chromosome and SHM occurs at the DNA level, it is interesting to observe that the mutational

382  tendencies are not reciprocal.

383     Having observed disproportional mutation tendencies in nucleotide pairs, we investigated the

384  mutability of reported motifs in a strand-specific manner (Material and methods). It is worth noting

385  that all nucleotides and motifs were extracted from the forward V gene sequences, and the reverse

386  sequences were discarded. Every nucleotide was classified exclusively in a single motif. Therefore,

387  the bases between categories have no overlap. We confirmed that SY<u>C</u> (where S is C or G; and Y is C

388  or T) and <u>G</u>RS (where R is A or G) are *bona fide* coldspots that showed the lowest frequency of

389  mutations. The motifs WR<u>C</u>Y, R<u>G</u>YW, W<u>A</u>, and <u>T</u>W also showed much higher mutations than did

390  coldspots as reported by others (Liu and Schatz, 2009; Pham et al., 2003). However, significantly

391  different mutabilities were observed again between reciprocal motifs (WR<u>C</u>Y and R<u>G</u>YW, W<u>A</u> and

392  <u>T</u>W) (Figure 6a, Figure S17c and Table S6). This result suggested strongly that SHM is introduced in

393  a strand-specific manner.

394     Complementarity-determining region (CDR) 1 and CDR2 exhibited higher mutations as expected

395  (Figure 6b and Figure S17d). While framework region (FR) 1 and FR2 displayed lower mutation rate

396  in contrast with CDRs, the region immediately adjacent to the CDR regions was also subjected to a

397  high frequency of mutations. These results support the idea that FRs provide the backbone of the

398  antibody, while CDRs accumulate mutations to achieve high affinity binding to a target antigen.

399  Consistent with previous observations (Shapiro et al., 1999), there was a considerable amount of SHM

400  in the FR3 region. Interestingly, the base with highest mutation rate was found near the end of the FR3.

401  A closer look at the germline sequence revealed that this nucleotide represents the third position

402  within a codon. The space for nucleotides, associated codons and amino acids was dominated by *G*,

403  *GTG*, and *V (valine)*, respectively (Figures S18a-c). Interestingly, however, in most cases, this site

404  does not occupy any previously identified canonical hotspot (Figure S18d). Nucleotide substitution

405  analysis at this locus showed no preference and has no impact on the encoded amino acid except for

406  eight alleles in the IGHV5 gene family (Figures S18e-g and S18h). In addition, the mutation spectrum

407    profiles varied in different IGHV families (Figure 6c).

408    Although the conservative substitution (transition within the same amino acid group) dominated

409    amino acid substitution profile, we observed relatively-high frequencies for non-conservative

410    substitutions, such as *H* to *Y* and *N* to *D*, that were not identified before using limited number of

411    datasets. In addition, we found *W* and *C* were least mutated (3.95% and 1.63%) and mutated to (1.00%

412    and 2.17%) (Figure 6c and Figure S17e).

413    The level of SHM as a function of confounding factors, such as age, sex, and isotypes, has also

414    been explored to some extent (Jiang et al., 2013; Kitaura et al., 2017; Wang et al., 2014). Nonetheless,

415    there have been no studies to date profiling SHM from a large data set. Using 363 samples from 290

416    donors, we reviewed the role of age, sex and isotypes on the frequency of SHM. We found negligible

417    differences between males and females, except for those between 41-50 years of age which probably

418    result from uncommon sampling bias. (Figure S19). We therefore combined data from male and

419    female donors to investigate the effects of age and isotype. Switched isotypes, namely IgG, IgA, and

420    IgE, had comparable level of SHM (7 - 8%), while IgM and IgD had much lower level of SHM (1 -

421    2%) (Figure S20), consistent with a previous report (Kitaura et al., 2017). We also observed a positive

422    correlation between the level of IgG SHM and age, except for individuals in the 41-50 age range

423    (Figure 6e). When we measure the contribution of age to SHM levels using a linear model, we

424    obtained r-square values of 0.37 and 0.28 for male and female, respectively. Despite the compromised

425    goodness of the model, we confirmed this correlation at the population level and estimated that the

426    SHM increases by approximately 0.05% each year. This increase would mean that in general, a parent

427    bears 1% more SHM than their children (Figure 6f).

428

# Discussion

430    Extremely large data sets have proven to be powerful for computational analysis to reveal

431    patterns, trends, and associations. The advent and application of high-throughput sequencing

432    technology advanced the study of complex biological systems, launching projects such as the 1000

433    Genomes Project, The Cancer Genome Atlas, the Encyclopedia of DNA Elements, and the NIH

434    Human Microbiome Project. Inspired by the success of these projects, we systematically analyzed the

435    largest antibody repertoire dataset to date and scrutinized, for the first time at this scale, the key

436    features of the antibody repertoire.

437    In addition to the uneven usage of germline genes, we identified a set of core V genes that

438    contribute to the clear majority of the repertoire. Although the other V genes are less frequently

439    observed in the current datasets, we believe their absence is the result of shallow sequencing depth

440    compared to the complexity of antibody repertoire. Nonetheless, these core and "rare" V gene sets

441    may serve as a reference for discovering gene usage fluctuations that are associated with or specific to

442    particular diseases. We found that the number of public clones between two repertoires also relied on

443    the sequencing depth. Moreover, a fraction of public clones identified in repertoire comparisons were

444    reported to be disease-associated or antigen-specific. This result supports the notion of antibody

445    convergence and also suggests that the antibody repertoire may help us trace an individual's immune

446    history and may therefore be useful in selecting vaccines and immunotherapy for certain diseases.

447    The fundamental B cell biology that underlies the specific patterns of germline usage, D-D fusion,

448    and SHMs revealed in our analyses remains controversial. Follow-up experiments may reveal the

449    mechanisms behind these phenomena and thus advance our understanding of B cell development as

450    well as its response to immune perturbations.

451    Due to the intrinsic amplification bias caused by different amplification strategies and various

452    primer sets, we did not perform analyses of clonal expansion, diversity, and evenness. A common

453    standard for both experimental design and bioinformatics analysis will be critical for future studies.

454    The human antibody repertoire possesses extreme diversity. Compared to the aforementioned

455    prior studies, the number of samples analyzed here is far from sufficient to capture all this diversity.

456    Moreover, antibody repertoires from a broad spectrum of diseases as well as different isotypes from

457    various tissue types are still needed for a better understanding of humoral immunity. Due to the

458 limited source of human samples, it is likely that further studies with model systems such as mouse,

459 rat, and macaque will bring us more insights.

## Acknowledgements

## Author Contributions

472 X. Y., Y. Z., H. Z., Y. Z., C. L., J. W., C. M., and Y. Z. performed the bioinformatics analyses on the

473 data. M. W., D. S., C. L., Y. D., S. G., L. X., R. W., and J. O. collected blood samples and conducted

474 the biological experiments. M. W., S. G., R. W., and C.J. L. prepared the libraries and performed

475 Illumina sequencing. C. C., W. Y., H. Z., J. C., L. Q., H. Z., J.X. B., L. W., G. C., X. Y. and Z. Z.

476 designed the project, biological experiments as well as bioinformatics analyses. X. Y., M. W., D. S., Y.

477 Z., H. Z., Y. Z., C. L., G. C., X. Y. and Z. Z. co-wrote the manuscripts.

## Declaration of Interests

479 The authors declare no competing financial interests. China Patents No. CN2019104688441,

480 CN2019104688441, and CN2019104688579.

# Figure titles and legends

481

482    **Figure 1. Overview of the enrolled datasets. (a)** The number of samples in each enrolled project.

483    The X axis shows NCBI SRA project IDs, and ZZHLAB indicates the antibody repertoires generated

484    in our lab. The Y axis shows the log10 transformed number of samples. The numbers of samples

485    excluded by data size and experimental design are shown in red and grey, respectively. **(b, c, d, e, f)**

486    show sample distribution based on **(b)** age; **(c)** sex; **(d)** PCR amplification strategy; **(e)** classification

487    (healthy or diseased); and **(f)** various tissue or blood.

488    **Figure 2. Germline gene usage and core V genes. (a)** The heatmaps show the normalized usage of V

489    (top panel), D (middle panel), and J (bottom panel) genes. Each column shows color-coded gene

490    usage for a dataset. Each row shows usage pattern of a particular gene (IDs labeled on the left side) in

491    different datasets. The bar graphs to the right of the heatmaps show the number of samples in which

492    each gene was present. J genes were present in almost every sample. The bar graph on top of the V

493    gene usage heatmap shows the number of V genes present in each sample. **(b)** and **(c)** The V genes

494    were ordered based on their occurrences in 582 individuals from high (left) to low (right). The X-axis

495    shows the number of high frequency V genes included. **(b)** The Y-axis shows the percent of total

496    clones that were represents by the most frequent V genes shown on X-axis. The color shows the log

497    10 transformed number of samples each pixel represents. **(c)** The red line indicates the median

498    fractions of total clones that were represents (left Y-axis) by the inclusion of top number of clones

499    shown on X-axis. The blue line represents the slope of median clone fraction variation (on the red line)

500    based on the adjacent 10 data points, 5 on the left and 5 on the right.

501    **Figure 3. Gene usage patterns with regard to genetic background, age, and gender. (a)** The

502    Pearson's correlation (Pearson's r) distribution of the gene usage between 5,261 paired samples. The

503    Pearson's r values were ordered from low to high. The red and light pink lines represent Pearson's r

504    values calculated using all V genes and 53 core genes, respectively. The blue and green dots indicate

505    the Pearson's r values between same and different cell types for monozygotic twins, respectively. **(b)**

506    The relationship between core V genes and ages. The X-axis shows V gene ordered by frequency

507    (Table S1). The Y-axis indicates the R2 values calculated for a particular V gene at different ages

508    (Supplementary. Fig. 6 and Materials and methods). **(c)**, **(d)** and **(e)** show comparisons of core V **(c)**,

509    D **(d)**, and J **(e)** genes between male and female. The red triangles indicate genes whose usage was

510    significantly different between sexes.

511    **Figure 4. Recombination and modification between V(D)J recombination. (a)** Recombination

512    count and frequency of different VD/DJ segments. The logarithm of the count is shown by the color of

513    the points, and the frequency of recombination is shown by the size of the points. The line at the

514    margin shows the number of each gene segment. V genes, core genes and non-core genes are marked.

515    The arrow shows the direction in IGH locus. **(b)** Distribution of CDR3 length in all sample, clones

516    with whole V, D, J assignment, and DD fusion. **(c)** Frequency of DD fusion in each isotype. The line

517    plot shows the number of samples with at least 5,000 clones in each isotype. **(d)** D gene usage in DD

518    fusion. **(e)** Frequency of DD fusion with different span; the span of adjacent D gene is 1. **(f)** The

519    number of DD fusions in all clones. The x axis represents the D gene at the 5' end, and the y axis

520    represents the D gene at the 3' end. The bar plot at the margin shows the number of each row or

521    column.

522    **Figure 5. Inter-sample abundance and gene usage of public clones. (a)** The heatmap in the center

523    indicates the abundance of public clones between samples. The top bar chart indicates the number of

524    recovered total clones for each sample. The number of public clones between each pair of sample has

525    been subjected to logarithmic transformation $(T=\log(1+Pab))$. The number of public clones between

526    samples within the same project has been set to 0 to remove chimera-related effects. Note that some

527    samples from PRJNA260985 and PRJNA280743, were predicted to come from the same donors and

528    the observed public clones between these samples was set to 0. **(b)** Linear model delineating the

529    correlation between inter-sample public clone abundance and the product of their clone abundance. **(c)**

530    Public clone size percentage as a function of donor sharing count.

531    **Figure 6. Somatic hypermutation patterns and influence factors. (a)** The stacked column diagram

532    shows the mutation percentage of motifs and composition of mutation targets. The X axis shows the

533    different motifs in germline sequences. The Y axis shows the composition of the mutated nucleotide

534    of this motif. The line chart shows the mutation fraction of every motif. The red-colored label

535    represents hot-spot, the blue colored label represents cold-spot. The underlined letter represents the

536    mutation site. **(b)** and **(c)** show the mutation rate among different functional alleles and families. **(b)**

537    The combined heatmap shows the mutation rate among used functional alleles in selected IgG samples.

538    Each column shows the position of completion of the V segment from FR1 to FR3. Each row shows

539    the functional alleles occurred in datasets. The area chart represents the average mutation rate in every

540    position. The bar graph left to the heatmap shows the family of occurred alleles which ordered by the

541    number of clones who were shows in the right bar graph. The color of the heat map represents the

542    mutation rate of every position from used functional alleles. **(c)** The X axis shows the position of the

543    V segment from FR1 to FR3. The Y axis shows the average mutation rate from different families. The

544    area chart shows the overall average mutation rate about used functional alleles. The red lines and blue

545    dotted lines show the result of the mutation rate of every family based on consensus and weight matrix

546    methods. **(d)** The combined heatmap shows the substitution among amino acid. Each column and each

547    row represents an amino acid. The germline residue is located on the x axis, and the mutated amino

548    acid is located on the Y axis. The line graph represents the ability of each amino acid to be mutated

549    and mutated. **(e)** The boxplot shows the mutation rate for different age groups across multiple

550    functional region and whole region. The points on top of each boxplot indicates the outliers. **(f)** The

551    scatter plot (orange for male and blue for female) shows the correlation between mutation rate and age.

552    Two lines in the figure are the predicted linear regression model for male and female. R-squared value

553    were marked on the top left in this figure.

554

# Materials and Methods

555

## Dataset enrollment criteria

556

557 We searched for bioprojects that were related to the antibody repertoire on the Sequence Read Archive

558 (SRA: https://www.ncbi.nlm.nih.gov/sra) from the National Center for Biotechnology Information

559 (NCBI: https://www.ncbi.nlm.nih.gov/). We identified thirty-eight projects before Feb 28, 2019. The

560 datasets from the included projects were subjected to two consecutive filter processes. The first filter

561 procedure was based entirely on sample metadata provided by SRA and the corresponding papers. The

562 criteria include:

563 ● Homo sapiens

564 ● Illumina platform

565 ● Pair-end Library Layout

566 ● Sequencing length >= 250

567 ● Natural sample directly extracted from human tissues (excludes those samples derived from cell

568 lines)

569 ● No specific amplification

570 ● Library source is either GENOMIC or TRANSCRIPTOME

571 ● No spike-in sequences

572 The second filter procedure was based on the results when preprocessing finished, criteria here

573 consists of,

574 ● Number of productive reads for heavy chain > =10,000

575 ● Fraction of heavy chain >= 20%

576 **In-house dataset**

577 **Subjects**

578 A total of 295 peripheral blood mononuclear cells (PBMCs) samples were collected. Of these, 254

579 were derived from healthy individuals (without recent infection events), 18 were from HBV patients,

580 16 were from H7N9 patients, 6 were from individuals involved in traffic accidents, and 1 was from a

581 patient with meningitis. Peripheral blood samples (1 ml) obtained from each volunteer were collected

582 in an EDTA-containing sterile tube and stored at room temperature for no more than 6 hours. PBMCs

583 were isolated by Ficoll-Paque density centrifugation using Lymphoprep™ (Axis-Shield, 1114547),

584 and the isolated cells were lysed in RLT buffer (Qiagen) supplemented with 1% β-mercaptoethanol

585 (Sigma) before being stored in -80□ for short-term storage. This protocol was approved by the Ethics

586 Committee at Southern Medical University. Informed consent was obtained from all participants.

587 **RNA extraction, reverse transcription, 5'RACE amplification, and next-generation sequencing**
588 **procedures**

589 RNA purification was carried out using the RNeasy Mini Kit (Qiagen, 74106) according to the

590 manufacturer's instructions. The concentration of the RNA was determined using a NanoDrop 2000c

591 Spectrophotometer (ThermoFisher Scientific). Five hundred nanograms of RNA purified from each

592 sample was used for cDNA synthesis with a total volume of 20 µl. cDNA was prepared using a

593 SMARTer RACE cDNA Amplification Kit (Clontech, 634928) according to the manufacturer's

594 instructions. Forward primers were synthesized according to SMARTer RACE protocol. The first 50

595 bp of the first constant domain (CH1) of heavy chain (IgG) were used to design the reverse primers.

596 We also designed 8-11bp barcode at the upstream of these primers to distinguish samples. One

597 microliter of the reverse transcription mixture was used as a template in a 20 µl PCR reaction. Primers

598 were used at a final concentration of 100 nM. The thermal cycling conditions were programmed as

599 follows: denaturation at 95°C for 3min, 30 cycles of denaturation at 98°C for 20s, annealing of primer

600 to DNA at 60°C for 15s, and extension by Kapa HiFi HotStart Ready Mix (KAPA Biosystems, kk2602)

601 at 72°C for 15s, followed by a final extension for 5 min at 72□. PCR products were analyzed on a 1.5%

602 agarose gel, and the appropriate bands (~600 bp) were purified using the Nucleospin Gel & PCR

603 Clean-up kit (Macherey-Nagel, 704609.25). DNA Concentration was measured using the NanoDrop

604 2000c Spectrophotometer (Thermo Fisher Scientific), and 400 ng of DNA was used to prepare

605 libraries using a Universal DNA Library Prep Kit for Illumina V3 (Vazyme, ND607-01), strictly

606 following the manufacturer's instructions. Libraries were quantified using the Qubit 4.0 fluorometer

607 (ThermoFisher Scientific) and re-quantified using the KAPA qPCR kit (KAPA Biosystems, 4824). The

608 size of adapter-ligated DNA fragments (approximately 800 bp) was determined using a Bioanalyzer

609 2100 system (Agilent). Each library was subjected to $2 \times 300$ bp paired-end sequencing using MiSeq

610 Reagent V3 kits (Illumina, MS-102-3003).

611 **Germline gene assignment and clonotype assemble**

612 Paired-end FASTQ files downloaded from SRA and generated by our laboratory were inputted into

613 *MiXCR* (version 3.0.7) and run with the following parameters:

614 Align: *mixcr align --library my_library -t 8 -r align_log.txt R1 R2 alignments.vdjca -s hs*

615 Assemble: *mixcr assemble -r assemble_log.txt -OseparateByV=true -OseparateByJ=true -Osepar*

616 *ateByC=true alignments.vdjca clones.clns*

617 Export clones: *mixcr exportClones clones.clns clones.txt*

618 Export Alignments: *mixcr exportAlignments -f -readIds -vHit -dHit -jHit -cHit -vGene -dGene*

619 *-jGene -nFeature CDR3 -aaFeature CDR3 -defaultAnchorPoints alignments.vdjca alignments.txt*

620 We built germline references for V, D, J, and C gene segments locally, and the germline refe

621 rences for V, D and J gene used in this study were customized using *repseqio* (v1.2.12, https:

622 //github.com/repseqio/repseqio). Reference sequences were obtained from IMGT/GENE DB (htt

623 p://www.imgt.org/genedb/) and are provided in Table S7. The formatted information for the re

624 ference constant region sequences was directly extracted from the *MiXCR* built-in reference (v

625 1.5) and then appended to the formatted customized reference for V, D and J genes. *MiXCR*

626 clustered sequences with the same V, J, and C allele assignment and CDR3 nucleotide sequen

627    ce into a clone with the parameters above. An in-house *Python* script was used to merge clo

628    nes with the same V and J gene and CDR3 nucleotide sequence into a clone. If we investiga

629    ted isotypes effect on some indices such gene usage, the C gene was also taken consideration.

630

### Comparison of gene usage and expression between Multiplex and RACE

632    Gene usage was defined as the number of clones with a given gene segment divided by the total

633    number of clones. Gene expression was defined as the number of reads with a gene divided by the

634    number of productive reads. For each gene segment, the median of usage and expression for either

635    Multiplex or RACE was used for linear regression. Usage or expression from Multiplex was defined

636    as independent variable while that from RACE was considered as dependent variable. The *regplot,*

637    *r2_score,* and *pearsonr* functions in *seaborn* (version 0.9.1) and the *sklearn* (version 0.20.2) and *scipy*

638    (version 1.2.1) *Python* modules were used to visualize the linear regression and to calculate R squared

639    values and a Pearson Correlation Coefficient.

### Overview and core gene set selection of gene usage

641    To show gene usage for all 2,152 samples clearly, we set thresholds for V, D, and J genes. If the usage

642    was greater than the threshold, we used the threshold value instead of the original value. The average

643    of the maximum of each sample for V, D, and J gene were calculated as thresholds. For core V gene

644    set selection, we first sorted genes according to their occurrence in 582 donors. We then enrolled 102

645    V genes one by one and computed the accumulated clone fraction with specific V genes for 2,152

646    samples. The median of clone fraction for all samples was selected and the slope of them was

647    computed. The slope of $x_i$ was equal to the distance of clone fraction at $x_{i-5}$ and $x_{i+5}$ divided by 11.

648    Finally, if the slope was less than 0.001, we determined the clone fraction arrived the plateau and

649    chose this gene set as the core gene.

650 **Features' effect on gene usage**

651 **Genetic background**

652 Two hundred and twenty-two healthy peripheral blood mononuclear cell samples, obtained from 29

653 male and 22 female individuals from 21 to 30 years of age, were amplified by RACE and used to

654 explore how genetic background affects gene usage. We examined two gene sets containing 53 core V

655 gene and 102 V genes. The Pearson Correlation Coefficient for every sample pair was calculated

656 using *pearsonr* from a Python module named *scipy* (version 1.2.1). Sample pairs from the same donor

657 were excluded. We performed statistics for male and female samples separately.

658 **Age**

659 We chose 499 healthy PBMC samples drawn from 94 male donors and 164 female donors by RACE.

660 Linear regression was conducted using *LinearRegression* and *r2_socre* from *sklearn* (version 0.20.2)

661 module. The independent and dependent variables were age and gene usage (of 53 V genes, 34 D

662 genes, and 6 J genes), respectively. Samples derived from males and females were analyzed

663 separately.

664 **Sex**

665 We performed two independent sample t-tests for 53 V genes, 34 D genes, and 6 J genes on the 499

666 samples selected above using *ttest_ind* from *scipy* (version 1.2.1). Genes whose P values were less

667 than 0.05 were defined as differentially used in the male samples and the female samples.

668 **Isotype**

669 We first obtained isotypes composition including IgA, IgD, IgE, IgM, and IgG for the 499 samples

670 above. Because the fraction of IgE was too low to compare, we discarded this isotype. Based on the

671 isotype fraction, there were 14 female and 37 male samples that could be used for this analysis. We

672 then merged clones for each isotype from different samples derived from the same donor and

673 recalculated gene usage for them. Gene usage was regarded as a vector, and *Euclidean distance* was

674 calculated using *DistanceMatrix* from *scikit-bio* (version 0.5.5) to measure the similarity of gene

675   usage between different isotypes. The *nj* function from *scikit-bio* was used to build a neighbor joining

676   tree. Finally, we used *Dendroscope* (version 3.6.3) to generate the trees. Samples analyzed for gene

677   usage related to genetic background, age, and so on, are shown in Table S8.

## VD/DJ recombination

679   To compare the recombination bias in all clones, we only analyzed clones with a D gene assignment.

680   Clones with full V, D, J gene assignments were pooled together. Clones without stop codons or out-of

681   –frame mutations in the CDR3 region were considered to be productive clones. If multiple

682   assignments occurred, the gene with the highest score was used for the analysis. The clones were

683   separated into subgroups according to VD/DJ recombination, and the frequency of each group was

684   calculated. The number of each gene was calculated according to about 352 million productive clones.

## D-D fusion detection

686   *IgScout* was used to detect D-D fusion. Input files were extracted from the *MiXCR* results using a

687   custom generated script, which used default parameters and the same reference as *MiXCR*. No other

688   filter was used to detect D-D fusion. The identity of D-D fusion and alignment length were calculated

689   by a custom script written in python. Levenshtein distance was used to quantify the difference

690   between the reference and the aligned sequences. The length of the aligned sequence was calculated

691   directly from the result of *IgScout*.

## Position bias in D-D fusion

693   Tandem CDR3s from all samples were pooled together to calculate gene usage in all D-D fusions. We

694   defined D1 as the D gene at the 5' end and D2 as the D gene at the 3' end. The span of two D genes

695   was defined as the genes position number between 5' D and 3'D on the corresponding chromosome.

696   The span of two adjacent D genes was 1. A negative value indicated that the D1 gene was located at

697   the 5' end of D2 on the chromosome, and a positive value represented a D2 gene located at the 3' end

698   of D1 on the chromosome.

**Comparison of D-D frequency between isotypes**

To compare the D-D frequency between isotypes, we also included C gene annotations. Due to the low frequency of D-D fusion, we only included samples that contained at least 5,000 clones for corresponding isotypes. Our analysis included 104, 594, 489, 163 samples for IgD, IgM, IgG, and IgA, respectively. IgE was not included in the analysis because none of the samples met our criteria of at least 5000 annotated clones. The frequency of D-D fusion in each sample was calculated as the number of D-D fusions in the corresponding isotype divided by the total number of corresponding isotypes.

**CDR3 length distribution**

The number of nucleotides in each clone was defined as the CDR3 length of the clone. The distribution of all clones was showed calculated from total clones. CDR3 clones with a D gene assignment were calculated as D-containing CDR3. The length of Tandem CDR3s was calculated from an output file named tandem_cdr3s.txt generated by *IgScout*. To make the distribution comparable, the frequency of CDR3s in each length group was used.

**Public clone abundance profile**

Public clones were those clonotypes (defined before) shared by at least two donors from two or more projects. Therefore, number of public clones between two samples from the same project or the same donor was set to zero. This strict criterion for clonotypes was applied to remove 'public' clones resulting from chimeric artifacts.

**Linear model to delineate the stochastic nature of gene recombination**

Linear models were constructed using only valid sample pairs that derived from different donors in different projects. Associated coefficients for regression equation and R squared were estimated using function *linear_model.LinearRegression* within the Python package *sklearn* (version 0.20.2).

### Profile of gene usage, CDR3 length, and Junctional Modification

Non-redundant public clones were used to profile gene usage and CDR3 length distribution, while redundant public clones were used for junctional modification analysis. The junctional modification calculation method is the same as above. Statistical analysis was carried out using a two-tailed unpaired Student's t-test.

### Antigen- or disease-related antibody database overlapping

We generated custom antigen- and disease-related antibody databases (unpublished results). All curated antibody sequences (heavy chain) were collected from following databases: i) IMGT/LIGM_DB (http://www.imgt.org/ligmdb/); ii) abYsis (http://www.bioinf.org.uk/abysis3.1/index.html); iii) EMBLIG (http://acrmwww.biochem.ucl.ac.uk/abs/abybank/emblig/); iv) bNAber (http://bnaber.org/), v) HIV_DB (http://www.hiv.lanl.gov/); vi) NCBI Nucleotide database (https://www.ncbi.nlm.nih.gov/nuccore); and vii) EBI ENA (https://www.ebi.ac.uk/ena). For now, it is comprised of 65,088 non-redundant antibody heavy chain sequences, corresponding to 53,579 unique CDR3 amino acid sequences and 163 types of antigen or disease, including HIV, hematological malignancies, preterm birth, influenza and etc. Antigen- or disease-related enrichment analysis of overlapping antibodies was performed using a hypergeometric model, implemented with the *stats.hypergeom.cdf* function within the Python package *scipy* (version 1.2.1). The false discovery rate was calculated using the *Benjamini-Hochberg* method implemented with an in-house script.

### Somatic hypermutation

#### Sample selection

Only samples which were from healthy donors' PBMC and were amplified using RACE protocol were included in the somatic hypermutation analysis. Since some experimentally qualified datasets

745     with which the alignment information failed to be exported, 363 samples were included in the end

746     (Table S8).

747     **Export alignment**

748     Alignment information used to measure somatic hypermutation was exported using *MiXCR* with the

749     following parameters:

750     Assemble: *mixcr assemble -r assemble_log.txt -OseparateByV=true -OseparateByJ=true*

751     *-OseparateByC=true -a alignments.vdjca clones.clna*

752     Export Alignments: *mixcr exportAlignments -f -readIds -cloneId -vHit -vAlignment -jHit -jAlignment*

753     *-cHit -cAlignment -nFeature FR1 -nFeature CDR1 -nFeature FR2 -nFeature CDR2 -nFeature FR3*

754     *-nFeature CDR3 -nFeature FR4 -aaFeature FR1 -aaFeature CDR1 -aaFeature FR2 -aaFeature*

755     *CDR2 -aaFeature FR3 -aaFeature CDR3 -aaFeature FR4 -defaultAnchorPoints clones.clna*

756     *alignments.txt*

757     **Quality filtering and data preprocessing:**

758     (1)  Read QC. Removal of reads that could not be merged by *MiXCR*, those without complete

759          variable regions (VR), those having been assigned with a pseudogene or a different V assignment

760          compared with their corresponding clones', those containing insertions or deletions and those

761          with stop codons or frameshifts in the variable region.

762     (2)  Clone QC. Removal of clones with only single qualified reads following read QC procedure

763          above.

764     (3)  VR deduplication. Deduplicating VR to obtain non-redundant sequence set

765     (4)  VR grouping. Grouping VR according to the isotypes reported in the clone files.

766 **Implementation of consensus and position-weighted matrix approaches**

767 The position-weighted matrix approach considers all qualified non-redundant reads within each clone.

768 Because each clone was a basic unit in the somatic hypermutation analysis, the mutation rate for a

769 certain position was calculated as the sum of mutation rate for all mutation events observed within

770 reads supporting this clone. The number of substitution types for a nucleotide (nt) or an amino acid (aa)

771 for a certain position was defined as 1 if a nucleotide or amino acid in a given position underwent the

772 same substitution event for all reads within a clone (with the same target nt or aa), otherwise it would

773 have a value less than 1.

774 For every clone, a theoretical consensus sequence was calculated based on the *motifs* module in

775 *Biopython* (version 1.73). The Hamming distance was used to calculate the distance between the

776 theoretical sequence and each true sequence, where the true sequence closest to the theoretical

777 sequence was taken as the representative sequence of the clone.

778 **Software**

779 In-house scripts were written in Python (version 3.7.4) based on the numpy (version 1.16.4),

780 Biopython (version 1.73), Levenshtein (version 0.12.0) and pandas (version 0.24.2) modules. To

781 visualize these results, we used the Python modules seaborn (version 0.9.1) and matplotlib (version

782 3.0.2) as well as GraphPad Prism (version 7.04).

# Supplemental Information titles and legends

784 **Figure S1. Age (a) and sex (b) composition of enrolled donors.** Total number of enrolled

785 donors is 582.

786 **Figure S2. The Pearson correlation coefficients of gene expression and usage between**

787 **Multiplex and RACE.** Left: the distribution of Pearson correlation coefficients of V (a), D (b),

788 and J (c) gene expression between Multiplex and RACE. Right: Pearson correlation coefficients

789    of V (d), D (e), and J (f) gene usage between these two groups. Note: There are 1,409 datasets

790    amplified by Multiplex and 743 datasets amplified by RACE.

791    **Figure S3. Fraction of repertoire containing the 53 core V genes.**

792    **Figure S4. Number of uncaptured V genes for 109 male (a) and 113 female (b) samples.**

793    **Figure S5. The Pearson correlation coefficients of D (a), J (b) for male and V (c), D (d), and**

794    **J (e) for female.** Lines in red show all genes, and the line in light purple shows the core V genes.

795    The dots indicate monozygotic twins from PRJNA300878. The blue dots indicate the same cell

796    type of monozygotic twins, and the green dots indicate different cell types (naïve and memory)

797    for them.

798    **Figure S6. Effect of Age on gene usage for females and males.** The R square of linear

799    regression between D (a), and J (b) gene usage and age in female. (c) The scatter plot for 53 V, 34

800    D, and 6 J genes usage and age. The X-axis means the age and the Y-axis stands for the usage. R

801    squared for V (d), D (e), and J (f) usage and age in the male.

802    **Figure S7. Gene usage in the infected and uninfected samples of the female (a, b, and c) and**

803    **male (d, e, and f).** Boxplots show usage from uninfected samples, while the red dots represent

804    usage from infected ones.

805    **Figure S8. The isotype composition of clone from the male (a) and the female (b).** Each

806    column shows one isotype including IgA, IgD, IgE, IgG, IgM, and None, and each row represents

807    a sample. For the row side color at the left of heatmap, the leftmost one was used to mark an

808    individual while the right one was used to distinguish a project. Note: None means those clones

809    cannot be aligned to a C gene.

810 **Figure S9. Clustering of V (a), D (b), and J (c) gene usage with different isotype in the male.**

811 IgA was labeled in light blue, IgG was filled in blue, IgD was colored in light green, and IgM was

812 labeled in green.

813 **Figure S10. Clustering of V (a), D (b), and J (c) gene usage with different isotype in female.**

814 IgA was labeled in light blue, IgG was filled in blue, IgD was colored in light green, and IgM was

815 labeled in green.

816 **Figure S11. Linear model for describing the correlation between number of public clones**

817 **and product of numbers of clones with each sample pair.** Only sample pairs with both clone

818 number being greater than (a) 10,000, (b) 100,000, (c)1,000,000, (d) 2,000,000, (e) 3,000,000, (f)

819 4,000,000 and (g) 5,000,000 were selected to demonstrate the correlation. The regression

820 functions are at the top of figures. Selected sample pair that with more clones show more fitness

821 of the linear model.

822 **Figure S12. Public clone gene segment usage.** (a, b) The two barplots show V and J gene usage

823 frequency between public and private clones. Gene segments have been sorted by overall

824 frequency and for v gene only those comprised more than 1% of the total repertoire were listed

825 here. (c, d) The scatter plot demonstrated gene usage frequency correlation between public and

826 private clones. The top left value ($\rho$) indicates Pearson's correlation coefficient.

827 **Figure S13. CDR3 nucleotide length distribution comparison between total productive**

828 **clones (n=267,761,654) and unique public clones (n=429,157).** Note that the upper limit for

829 length is determined according to a threshold of 1%. Two-sample Kolmogorov-Smirnov tests

830 were performed to investigate length distribution difference (P-value <1.149e-13).

831 **Figure S14. Junctional modification comparison between total clones from 2,165 samples**

832 **and public clones.** The boxplot in each subfigure demonstrates the distribution of each kind of

833   junction modification length, as indicated by the schematic in bottom right (a). (b) Non-templated

834   insertion length. (c) Palindromic insertion length distributions. (d) Deletion length distributions.

835   **Figure S15. Percent of clones with D hit(s) for 2,165 samples and all public clones.** The

836   boxplot on the top demonstrates the percent distribution for 2,165 samples (with a median of

837   97.9%), and the red point at the bottom indicates the percent for public clones (67.6%).

838   **Figure S16. Antigen-specific or disease-associated annotation of public clones.** (a)

839   Overlapping of unique clonotypes between public clones and antibody sequences curated with

840   related antigen or disease information. A clonotype here was defined as a unique CDR3 amino

841   acid sequence with deprecated conserved residuals at both ends. (b) Disease or antigen percentage

842   of annotated public clones. Terms in the same line in the legend are indicated by same color.

843   Terms in legend match the pie chart from top to bottom and from left to right.

844   **Figure S17. Somatic hypermutation patterns and influence factors.** (a) and (b) represent the

845   transform among nucleotides based on the algorithm of consensus and position weight matrix

846   (PWM). Each column and each row represent a nucleotide. The germline nucleotide is located on

847   the x axis and the mutated nucleotide is located on the Y axis. The total mutation rate and target

848   preference of every nucleotide are marked in the figure. (c), (d) and (e) used position weight

849   matrix (PWM) to describe the patterns of somatic hypermutation.

850   **Figure S18. Description of the 290th position, which has the highest mutation rate.** (a), (b),

851   (c), and (d) show the composition of nucleotide, amino acid, codon, and motif in the germline

852   sequence sorted by ratio, respectively. (e), (f) and (g) showed the fraction and composition of

853   mutation from different families. (h) The boxplot shows the comparison of mutation rate between

854   synonymous mutations and nonsynonymous mutations from different families.

855   **Figure S19. Mutation rate comparison between male and female based on IgG clones.** (a-e)

856   were based on consensus approach and (f-k) were based on PWM approach. Comparison were

857 performed independently in different age groups to remove age-related effect (marked on the top

858 right of each subfigure). and only those age groups with at least 10 donors for both genders were

859 presented here. The four numbers following each isotype in figure legends represent the number

860 of clones, samples, donors and projects, respectively.

861 **Figure S20. Mutation rate comparison between different isotypes.** (a-c) were based on

862 consensus approach and (d-f) were based on PWM approach. Comparison were performed

863 independently in different age groups to remove age-related effect (marked on the top right of

864 each subfigure). Only 3 age groups have two or more kinds of clones. The four numbers

865 following each isotype show the number of clones, samples, donors and projects, respectively.

866

867 # References

868 Bonsignori, M., Zhou, T., Sheng, Z., Chen, L., Gao, F., Joyce, M.G., Ozorowski, G., Chuang, G.Y.,

869 Schramm, C.A., and Wiehe, K.*, et al.* (2016). Maturation Pathway from Germline to Broad HIV-1

870 Neutralizer of a CD4-Mimic Antibody. CELL *165*, 449-463.

871 Boyd, S.D., Gaeta, B.A., Jackson, K.J., Fire, A.Z., Marshall, E.L., Merker, J.D., Maniar, J.M., Zhang,

872 L.N., Sahaf, B., and Jones, C.D.*, et al.* (2010). Individual variation in the germline Ig gene repertoire

873 inferred from variable region gene rearrangements. J IMMUNOL *184*, 6986-6992.

874 Briney, B., Inderbitzin, A., Joyce, C., and Burton, D.R. (2019). Commonality despite exceptional

875 diversity in the baseline human antibody repertoire. NATURE *566*, 393-397.

876 Briney, B.S., Willis, J.R., Hicar, M.D., Thomas, J.W., and Crowe, J.E. (2012). Frequency and genetic

877 characterization of V(DD)J recombinants in the human peripheral blood antibody repertoire.

878 IMMUNOLOGY *137*, 56-64.

879 Briney, B.S., Willis, J.R., McKinney, B.A., and Crowe, J.J. (2012). High-throughput antibody

880 sequencing reveals genetic evidence of global regulation of the naive and memory repertoires that

881 extends across individuals. GENES IMMUN *13*, 469-473.

882 Bürckert, J., Dubois, A.R.S.X., Faison, W.J., Farinelle, S., Charpentier, E., Sinner, R.,

883 Wienecke-Baldacchino, A., and Muller, C.P. (2017). Functionally Convergent B Cell Receptor

884 Sequences in Transgenic Rats Expressing a Human B Cell Repertoire in Response to Tetanus Toxoid

885 and Measles Antigens. FRONT IMMUNOL *8*.

886 Chothia, C., Lesk, A.M., Tramontano, A., Levitt, M., Smith-Gill, S.J., Air, G., Sheriff, S., Padlan,

887 E.A., Davies, D., and Tulip, W.R.*, et al.* (1989). Conformations of immunoglobulin hypervariable

888 regions. NATURE *342*, 877-883.

889 Early, P., Huang, H., Davis, M., Calame, K., and Hood, L. (1980). An immunoglobulin heavy chain

890 variable region gene is generated from three segments of DNA: VH, D and JH. CELL *19*, 981-992.

891 Faham, M., Zheng, J., Moorhead, M., Carlton, V.E., Stow, P., Coustan-Smith, E., Pui, C.H., and

892 Campana, D. (2012). Deep-sequencing approach for minimal residual disease detection in acute

893 lymphoblastic leukemia. BLOOD *120*, 5173-5180.

894 Gawad, C., Pepin, F., Carlton, V.E.H., Klinger, M., Logan, A.C., Miklos, D.B., Faham, M., Dahl, G.,

895 and Lacayo, N. (2012). Massive evolution of the immunoglobulin heavy chain locus in children with

896 B precursor acute lymphoblastic leukemia. BLOOD *120*, 4407-4417.

897 Glanville, J., Kuo, T.C., von Budingen, H.C., Guey, L., Berka, J., Sundar, P.D., Huerta, G., Mehta,

898 G.R., Oksenberg, J.R., and Hauser, S.L.*, et al.* (2011). Naive antibody gene-segment frequencies are

899 heritable and unaltered by chronic lymphocyte ablation. Proceedings of the National Academy of

900 Sciences *108*, 20066-20071.

901 Greiff, V., Miho, E., Menzel, U., and Reddy, S.T. (2015). Bioinformatic and Statistical Analysis of

902 Adaptive Immune Repertoires. TRENDS IMMUNOL *36*, 738-749.

903 Greiff, V., Weber, C.R., Palme, J., Bodenhofer, U., Miho, E., Menzel, U., and Reddy, S.T. (2017).

904 Learning the High-Dimensional Immunogenomic Features That Predict Public and Private Antibody

905 Repertoires. J IMMUNOL *199*, 2985-2997.

906 Hansen, T.O., Lange, A.B., and Barington, T. (2015). Sterile DJH rearrangements reveal that distance

907 between gene segments on the human Ig H chain locus influences their ability to rearrange. J

908 IMMUNOL *194*, 973-982.

909    He, L., Sok, D., Azadnia, P., Hsueh, J., Landais, E., Simek, M., Koff, W.C., Poignard, P., Burton,

910    D.R., and Zhu, J. (2015). Toward a more accurate view of human B-cell repertoire by next-generation

911    sequencing, unbiased repertoire capture and single-molecule barcoding. SCI REP-UK *4*.

912    Hoh, R.A., Joshi, S.A., Liu, Y., Wang, C., Roskin, K.M., Lee, J., Pham, T., Looney, T.J., Jackson,

913    K.J.L., and Dixit, V.P.*, et al.* (2016). Single B-cell deconvolution of peanut-specific antibody

914    responses in allergic patients. J ALLERGY CLIN IMMUN *137*, 157-167.

915    Hong, B., Wu, Y., Li, W., Wang, X., Wen, Y., Jiang, S., Dimitrov, D.S., and Ying, T. (2018).

916    In-Depth Analysis of Human Neonatal and Adult IgM Antibody Repertoires. FRONT IMMUNOL *9*,

917    128.

918    Jackson, K.J., Kidd, M.J., Wang, Y., and Collins, A.M. (2013). The shape of the lymphocyte receptor

919    repertoire: lessons from the B cell receptor. FRONT IMMUNOL *4*, 263.

920    Jackson, K.J., Liu, Y., Roskin, K.M., Glanville, J., Hoh, R.A., Seo, K., Marshall, E.L., Gurley, T.C.,

921    Moody, M.A., and Haynes, B.F.*, et al.* (2014). Human responses to influenza vaccination show

922    seroconversion signatures and convergent antibody rearrangements. CELL HOST MICROBE *16*,

923    105-114.

924    Jiang, N., He, J., Weinstein, J.A., Penland, L., Sasaki, S., He, X.S., Dekker, C.L., Zheng, N.Y., Huang,

925    M., and Sullivan, M.*, et al.* (2013). Lineage structure of the human antibody repertoire in response to

926    influenza vaccination. SCI TRANSL MED *5*, 119r-171r.

927    Joyce, M.G., Wheatley, A.K., Thomas, P.V., Chuang, G., Soto, C., Bailer, R.T., Druz, A., Georgiev,

928    I.S., Gillespie, R.A., and Kanekiyo, M.*, et al.* (2016). Vaccine-Induced Antibodies that Neutralize

929    Group 1 and Group 2 Influenza A Viruses. CELL *166*, 609-623.

930    Kidd, M.J., Jackson, K.J., Boyd, S.D., and Collins, A.M. (2016). DJ Pairing during VDJ

931    Recombination Shows Positional Biases That Vary among Individuals with Differing IGHD Locus

932    Immunogenotypes. J IMMUNOL *196*, 1158-1164.

933    Kitaura, K., Yamashita, H., Ayabe, H., Shini, T., Matsutani, T., and Suzuki, R. (2017). Different

934    Somatic Hypermutation Levels among Antibody Subclasses Disclosed by a New Next-Generation

935    Sequencing-Based Antibody Repertoire Analysis. FRONT IMMUNOL *8*.

936     Krebs, S.J., Kwon, Y.D., Schramm, C.A., Law, W.H., Donofrio, G., Zhou, K.H., Gift, S., Dussupt, V.,

937     Georgiev, I.S., and Schätzle, S.*, et al.* (2019). Longitudinal Analysis Reveals Early Development of

938     Three MPER-Directed Neutralizing Antibody Lineages from an HIV-1-Infected Individual.

939     IMMUNITY *50*, 677-691.

940     Kurtz, D.M., Green, M.R., Bratman, S.V., Scherer, F., Liu, C.L., Kunder, C.A., Takahashi, K., Glover,

941     C., Keane, C., and Kihira, S.*, et al.* (2015). Noninvasive monitoring of diffuse large B-cell lymphoma

942     by immunoglobulin high-throughput sequencing. BLOOD *125*, 3679-3687.

943     Larimore, K., McCormick, M.W., Robins, H.S., and Greenberg, P.D. (2012). Shaping of human

944     germline IgH repertoires revealed by deep sequencing. J IMMUNOL *189*, 3221-3230.

945     Laserson, U., Vigneault, F., Gadala-Maria, D., Yaari, G., Uduman, M., Vander, H.J., Kelton, W., Taek,

946     J.S., Liu, Y., and Laserson, J.*, et al.* (2014). High-resolution antibody dynamics of vaccine-induced

947     immune responses. Proc Natl Acad Sci U S A *111*, 4928-4933.

948     Li, G.M., Chiu, C., Wrammert, J., McCausland, M., Andrews, S.F., Zheng, N.Y., Lee, J.H., Huang,

949     M., Qu, X., and Edupuganti, S.*, et al.* (2012). Pandemic H1N1 influenza vaccine induces a recall

950     response in humans that favors broadly cross-reactive memory B cells. Proceedings of the National

951     Academy of Sciences *109*, 9047-9052.

952     Liu, M., and Schatz, D.G. (2009). Balancing AID and DNA repair during somatic hypermutation.

953     TRENDS IMMUNOL *30*, 173-181.

954     Liu, X., Zhang, W., Zeng, X., Zhang, R., Du, Y., Hong, X., Cao, H., Su, Z., Wang, C., and Wu, J.*, et*

955     *al.* (2016). Systematic Comparative Evaluation of Methods for Investigating the TCRβ Repertoire.

956     PLOS ONE *11*, e152464.

957     Madi, A., Shifrut, E., Reich-Zeliger, S., Gal, H., Best, K., Ndifon, W., Chain, B., Cohen, I.R., and

958     Friedman, N. (2014). T-cell receptor repertoires share a restricted set of public and abundant CDR3

959     sequences that are associated with self-related immunity. GENOME RES *24*, 1603-1612.

960   Maecker, H.T., Lindstrom, T.M., Robinson, W.H., Utz, P.J., Hale, M., Boyd, S.D., Shen-Orr, S.S., and

961   Fathman, C.G. (2012). New tools for classification and monitoring of autoimmune diseases. NAT

962   REV RHEUMATOL *8*, 317-328.

963   Miho, E., Roskar, R., Greiff, V., and Reddy, S.T. (2019). Large-scale network analysis reveals the

964   sequence space architecture of antibody   repertoires. NAT COMMUN *10*, 1321.

965   Parameswaran, P., Liu, Y., Roskin, K.M., Jackson, K.K.L., Dixit, V.P., Lee, J., Artiles, K.L., Zompi,

966   S., Vargas, M.J., and Simen, B.B.*, et al.* (2013). Convergent Antibody Signatures in Human Dengue.

967   CELL HOST MICROBE *13*, 691-700.

968   Patil, S.U., Ogunniyi, A.O., Calatroni, A., Tadigotla, V.R., Ruiter, B., Ma, A., Moon, J., Love, J.C.,

969   and Shreffler, W.G. (2015). Peanut oral immunotherapy transiently expands circulating Ara h 2‒

970   specific B cells with a homologous repertoire in unrelated subjects. J ALLERGY CLIN IMMUN *136*,

971   125-134.

972   Pham, P., Bransteitter, R., Petruska, J., and Goodman, M.F. (2003). Processive AID-catalysed

973   cytosine deamination on single-stranded DNA simulates somatic hypermutation. NATURE *424*,

974   103-107.

975   Pugh-Bernard, A.E., Silverman, G.J., Cappione, A.J., Villano, M.E., Ryan, D.H., Insel, R.A., and Sanz,

976   I. (2001). Regulation of inherently autoreactive VH4-34 B cells in the maintenance of human   B cell

977   tolerance. J CLIN INVEST *108*, 1061-1070.

978   Reddy, S.T., Ge, X., Miklos, A.E., Hughes, R.A., Kang, S.H., Hoi, K.H., Chrysostomou, C.,

979   Hunicke-Smith, S.P., Iverson, B.L., and Tucker, P.W.*, et al.* (2010). Monoclonal antibodies isolated

980   without screening by analyzing the variable-gene repertoire of plasma cells. NAT BIOTECHNOL *28*,

981   965-969.

982   Robins, H. (2013). Immunosequencing: applications of immune repertoire deep sequencing. CURR

983   OPIN IMMUNOL *25*, 646-652.

984   Saada, R., Weinberger, M., Shahaf, G., and Mehr, R. (2007). Models for antigen receptor gene

985   rearrangement: CDR3 length. IMMUNOL CELL BIOL *85*, 323-332.

986  Safonova, Y., and Pevzner, P.A. (2019). De novo Inference of Diversity Genes and Analysis of

987  Non-canonical V(DD)J Recombination in Immunoglobulins. FRONT IMMUNOL *10*, 987.

988  Schramm, C.A., and Douek, D.C. (2018). Beyond Hot Spots: Biases in Antibody Somatic

989  Hypermutation and Implications for Vaccine Design. FRONT IMMUNOL *9*, 1876.

990  Setliff, I., McDonnell, W.J., Raju, N., Bombardi, R.G., Murji, A.A., Scheepers, C., Ziki, R., Mynhardt,

991  C., Shepherd, B.E., and Mamchak, A.A*., et al.* (2018). Multi-Donor Longitudinal Antibody Repertoire

992  Sequencing Reveals the Existence of Public Antibody Clonotypes in HIV-1 Infection. CELL HOST

993  MICROBE *23*, 845-854.

994  Shapiro, G.S., Aviszus, K., Ikle, D., and Wysocki, L.J. (1999). Predicting regional mutability in

995  antibody V genes based solely on di- and trinucleotide sequence composition. J IMMUNOL *163*,

996  259-268.

997  Soto, C., Bombardi, R.G., Branchizio, A., Kose, N., Matta, P., Sevy, A.M., Sinkovits, R.S., Gilchuk,

998  P., Finn, J.A., and Crowe, J.E. (2019). High frequency of shared clonotypes in human B cell receptor

999  repertoires. NATURE.

1000  Souto-Carneiro, M.M., Sims, G.P., Girschik, H., Lee, J., and Lipsky, P.E. (2005). Developmental

1001  Changes in the Human Heavy Chain CDR3. The Journal of Immunology *175*, 7425-7436.

1002  Stern, J.N.H., Yaari, G., Vander Heiden, J., Church, G., Donahue, W.F., Hintzen, R., Huttner, A.J.,

1003  Laman, J., Nagra, R.M., and Nylander, A*., et al.* (2014). B cells populating the multiple sclerosis brain

1004  mature in the draining cervical lymph nodes. SCI TRANSL MED *6*, 107r-248r.

1005  Sui, J., Hwang, W.C., Perez, S., Wei, G., Aird, D., Chen, L., Santelli, E., Stec, B., Cadwell, G., and

1006  Ali, M*., et al.* (2009). Structural and functional bases for broad-spectrum neutralization of avian and

1007  human influenza A viruses. NAT STRUCT MOL BIOL *16*, 265-273.

1008  Tipton, C.M., Fucile, C.F., Darce, J., Chida, A., Ichikawa, T., Gregoretti, I., Schieferl, S., Hom, J.,

1009  Jenks, S., and Feldman, R.J*., et al.* (2015). Diversity, cellular origin and autoreactivity of

1010  antibody-secreting cell population expansions in acute systemic lupus erythematosus. NAT

1011  IMMUNOL *16*, 755-765.

1012  Tonegawa, S. (1983). Somatic generation of antibody diversity. NATURE *302*, 575-581.

1013 Truck, J., Ramasamy, M.N., Galson, J.D., Rance, R., Parkhill, J., Lunter, G., Pollard, A.J., and Kelly,

1014 D.F. (2015). Identification of antigen-specific B cell receptor sequences using public repertoire

1015 analysis. J IMMUNOL *194*, 252-261.

1016 von Büdingen, H., Kuo, T.C., Sirota, M., van Belle, C.J., Apeltsin, L., Glanville, J., Cree, B.A.,

1017 Gourraud, P., Schwartzburg, A., and Huerta, G.*, et al.* (2012). B cell exchange across the blood-brain

1018 barrier in multiple sclerosis. The Journal of clinical investigation *122*, 4533-4543.

1019 Wang, C., Liu, Y., Xu, L.T., Jackson, K.J., Roskin, K.M., Pham, T.D., Laserson, J., Marshall, E.L.,

1020 Seo, K., and Lee, J.Y.*, et al.* (2014). Effects of aging, cytomegalovirus infection, and EBV infection

1021 on human B cell repertoires. J IMMUNOL *192*, 603-611.

1022 Weinstein, J.A., Jiang, N., White, R.A., Fisher, D.S., and Quake, S.R. (2009). High-Throughput

1023 Sequencing of the Zebrafish Antibody Repertoire. SCIENCE *324*, 807-810.

1024 Wu, X., Zhang, Z., Schramm, C.A., Joyce, M.G., Do Kwon, Y., Zhou, T., Sheng, Z., Zhang, B., O

1025 Dell, S., and McKee, K.*, et al.* (2015). Maturation and Diversity of the VRC01-Antibody Lineage over

1026 15 Years of Chronic HIV-1 Infection. CELL *161*, 470-485.

1027 Wu, X., Zhou, T., Zhu, J., Zhang, B., Georgiev, I., Wang, C., Chen, X., Longo, N.S., Louder, M., and

1028 McKee, K.*, et al.* (2011a). Focused Evolution of HIV-1 Neutralizing Antibodies Revealed by

1029 Structures and Deep Sequencing. SCIENCE *333*, 1593-1602.

1030 Wu, X., Zhou, T., Zhu, J., Zhang, B., Georgiev, I., Wang, C., Chen, X., Longo, N.S., Louder, M., and

1031 McKee, K.*, et al.* (2011b). Focused Evolution of HIV-1 Neutralizing Antibodies Revealed by

1032 Structures and Deep Sequencing. SCIENCE *333*, 1593-1602.

1033 Wu, Y.B., James, L.K., Vander Heiden, J.A., Uduman, M., Durham, S.R., Kleinstein, S.H., Kipling,

1034 D., and Gould, H.J. (2014). Influence of seasonal exposure to grass pollen on local and peripheral

1035 blood IgE repertoires in patients with allergic rhinitis. J ALLERGY CLIN IMMUN *134*, 604-612.

1036 Zhu, J., Wu, X., Zhang, B., McKee, K., O'Dell, S., Soto, C., Zhou, T., Casazza, J.P., Mullikin, J.C.,

1037 and Kwong, P.D.*, et al.* (2013a). De novo identification of VRC01 class HIV-1-neutralizing

1038 antibodies by next-generation sequencing of B-cell transcripts. Proceedings of the National Academy

1039 of Sciences *110*, E4088-E4097.

1040   Zhu, J., Wu, X., Zhang, B., McKee, K., O'Dell, S., Soto, C., Zhou, T., Casazza, J.P., Mullikin, J.C.,

1041   and Kwong, P.D.*, et al.* (2013b). De novo identification of VRC01 class HIV-1-neutralizing

1042   antibodies by next-generation sequencing of B-cell transcripts. Proceedings of the National Academy

1043   of Sciences *110*, E4088-E4097.

1044

1045
1046   **DATA AVAILABILITY**

1047   In-house generated datasets are available at the NCBI Sequencing Read Archive

1048   (www.ncbi.nlm.nih.gov/sra) under BioProject number PRJNA564936. A table linking dataset

1049   accessions to their corresponding sample ids was provided in Table S9.

1050

# Main figures

Figure 1



**Figure 1 Overview of the enrolled datasets.** (a) The number of samples in each enrolled project. The X axis shows NCBI SRA project IDs, and ZZHLAB indicates the antibody repertoires generated in our lab. The Y axis shows the log10 transformed number of samples. The numbers of samples excluded by data size and experimental design are shown in red and grey, respectively. (b, c, d, e, f) show sample distribution based on (b) age; (c) sex; (d) PCR amplification strategy; (e) classification (healthy or diseased); and (f) various tissue or blood.

Figure 2



**Figure 2 Germline gene usage and core V genes. (a)** The heatmaps show the normalized usage of V (top panel), D (middle panel), and J (bottom panel) genes. Each column shows color-coded gene usage for a dataset. Each row shows usage pattern of a particular gene (IDs labeled on the left side) in different datasets. The bar graphs to the right of the heatmaps show the number of samples in which each gene was present. J genes were present in almost every sample. The bar graph on top of the V gene usage heatmap shows the number of V genes present in each sample. **(b)** and **(c)** The V genes were ordered based on their occurrences in 582 individuals from high (left) to low (right). The X-axis shows the number of high frequency V genes included. **(b)** The Y-axis shows the percent of total clones that were represents by the most frequent V genes shown on X-axis. The color shows the log 10 transformed number of samples each pixel represents. **(c)** The red line indicates the median fractions of total clones that were represents (left Y-axis) by the inclusion of top number of clones shown on X-axis. The blue line represents the slope of median clone fraction variation (on the red line) based on the adjacent 10 data points, 5 on the left and 5 on the right.
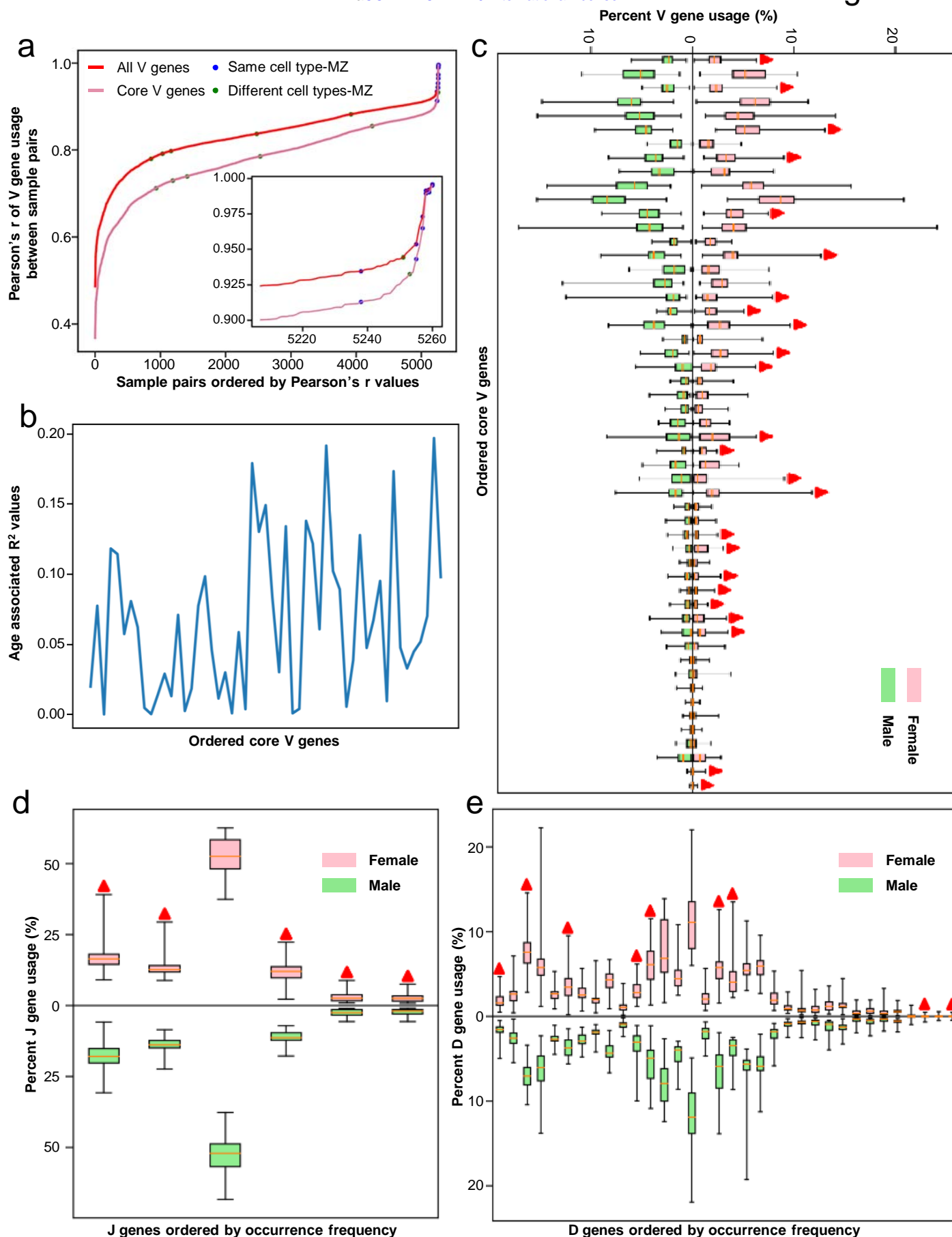
# Figure 3



**Figure 3 Gene usage patterns with regard to genetic background, age, and gender. (a)** The Pearson's correlation (Pearson's r) distribution of the gene usage between 5,261 paired samples. The Pearson's r values were ordered from low to high. The red and light pink lines represent Pearson's r values calculated using all V genes and 53 core genes, respectively. The blue and green dots indicate the Pearson's r values between same and different cell types for monozygotic twins, respectively. **(b)** The relationship between core V genes and ages. The X-axis shows V gene ordered by frequency (Table 1). The Y-axis indicates the $R^2$ values calculated for a particular V gene at different ages (Supp. Fig. 6 and Materials and methods). **(c)**, **(d)**, and **(e)** show comparisons of core V **(c)**, D **(d)**, and J **(e)** genes between male and female. The red triangles indicate genes whose usage was significantly different between sexes.
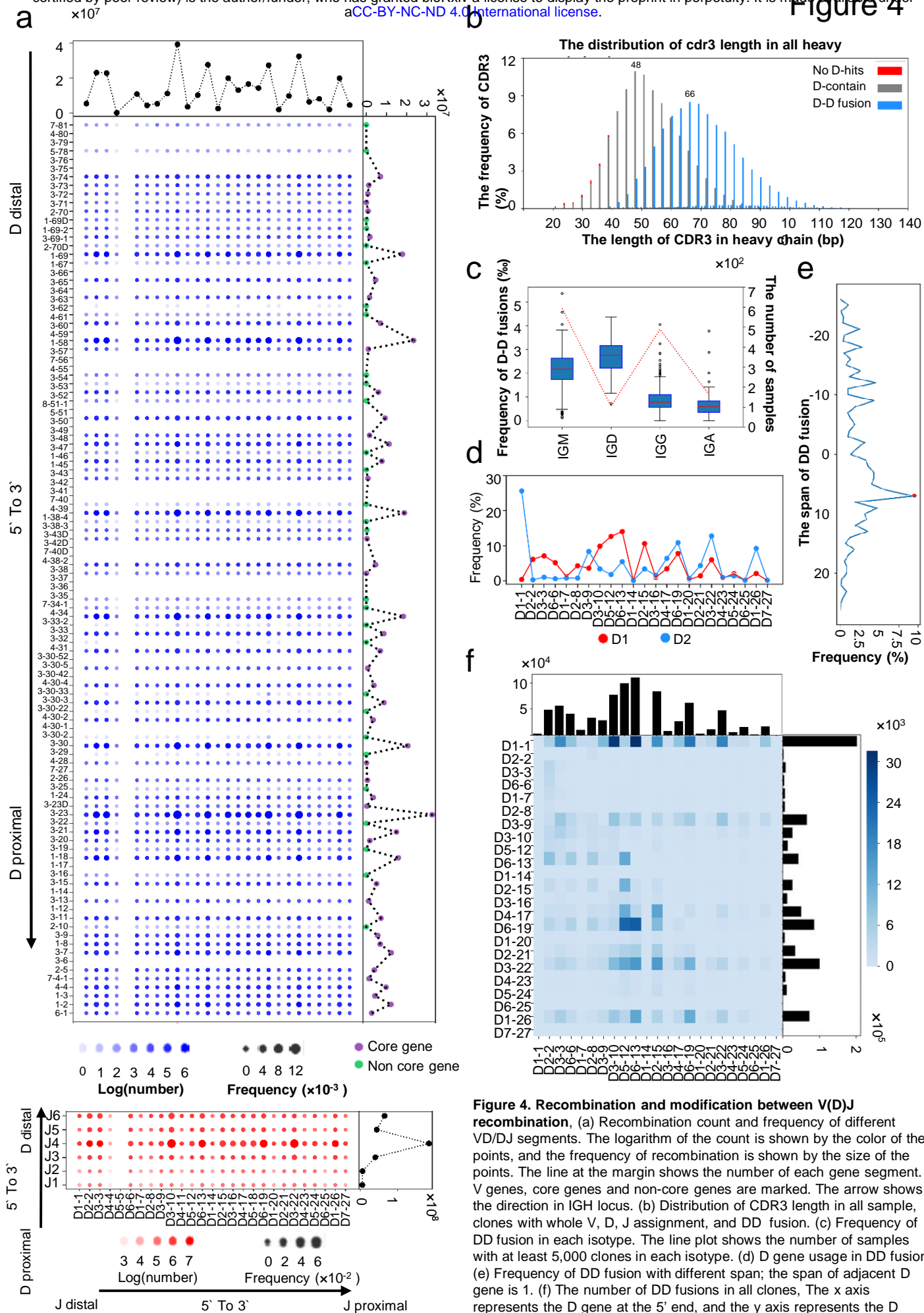
Figure 4



**Figure 4. Recombination and modification between V(D)J recombination**, (a) Recombination count and frequency of different VD/DJ segments. The logarithm of the count is shown by the color of the points, and the frequency of recombination is shown by the size of the points. The line at the margin shows the number of each gene segment. V genes, core genes and non-core genes are marked. The arrow shows the direction in IGH locus. (b) Distribution of CDR3 length in all sample, clones with whole V, D, J assignment, and DD fusion. (c) Frequency of DD fusion in each isotype. The line plot shows the number of samples with at least 5,000 clones in each isotype. (d) D gene usage in DD fusion. (e) Frequency of DD fusion with different span; the span of adjacent D gene is 1. (f) The number of DD fusions in all clones, The x axis represents the D gene at the 5' end, and the y axis represents the D gene at the 3' end. The bar plot at the margin shows the number of each
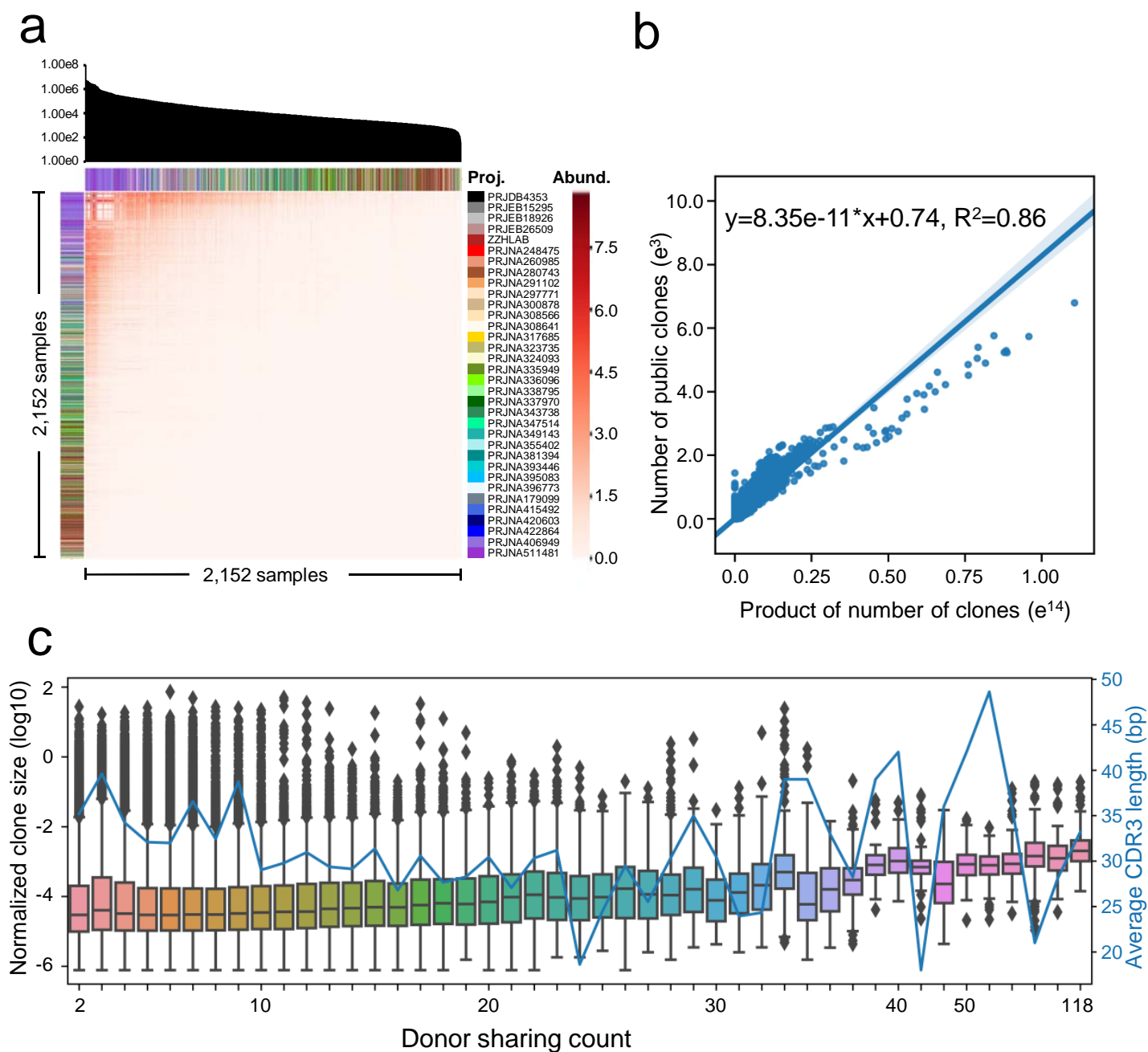
Figure 5



**Figure 5. Inter-sample abundance and gene usage of public clones. (a)** The heatmap in the center indicates the abundance of public clones between samples. The top bar chart indicates the number of recovered total clones for each sample. The number of public clones between each pair of sample has been subjected to logarithmic transformation (T=log(1+Pab)). The number of public clones between samples within the same project has been set to 0 to remove chimera-related effects. Note that some samples from PRJNA260985 and PRJNA280743, were predicted to come from the same donors and the observed public clones between these samples was set to 0. **(b)** Linear model delineating the correlation between inter-sample public clone abundance and the product of their clone abundance. **(c)** Public clone size percentage as a function of donor sharing count.
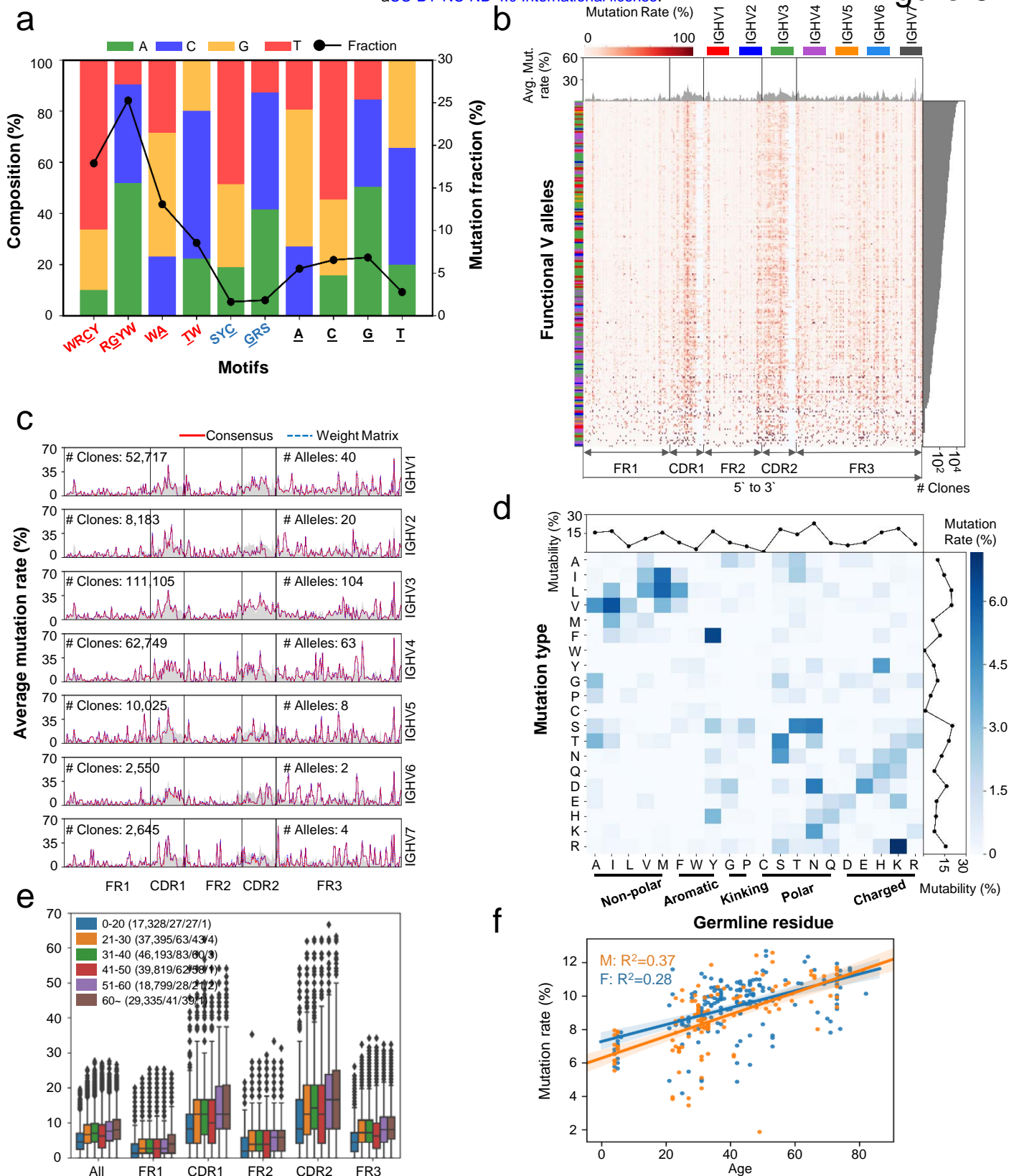
# Figure 6



**Figure 6 | Somatic hypermutation patterns and influence factors. (a)** The stacked column diagram shows the mutation percentage of motifs and composition of mutation targets. The X axis shows the different motifs in germline sequences. The Y axis shows the composition of the mutated nucleotide of this motif. The line chart shows the mutation fraction of every motif. The red-colored label represents hot-spot, the blue colored label represents cold-spot. The underlined letter represents the mutation site. **(b)** and **(c)** show the mutation rate among different functional alleles and families. **(b)** The combined heatmap shows the mutation rate among used functional alleles in selected IgG samples. Each column shows the position of completion of the V segment from FR1 to FR3. Each row shows the functional alleles occurred in datasets. The area chart represents the average mutation rate in every position. The bar graph left to the heatmap shows the family of occurred alleles which ordered by the number of clones who were shows in the right bar graph. The color of the heat map represents the mutation rate of every position from used functional alleles. **(c)** The X axis shows the position of the V segment from FR1 to FR3. The Y axis shows the average mutation rate from different families. The area chart shows the overall average mutation rate about used functional alleles. The red lines and blue dotted lines show the result of the mutation rate of every family based on consensus and weight matrix methods. **(d)** The combined heatmap shows the substitution among amino acid. Each column and each row represents an amino acid. The germline residue is located on the x axis, and the mutated amino acid is located on the Y axis. The line graphs represents the ability of each amino acid to be mutated and mutated. **(e)**. The boxplot shows the mutation rate for different age groups across multiple functional region and whole region. The points on top of each boxplot indicates the outliers. (f) The scatter plot (orange for male and blue for female) shows the correlation between mutation rate and age. Two lines in the figure are the predicted linear regression model for male and female. R-squared value were marked on the top left in this figure.