### 1 SARS-Cov-2-, HIV-1-, Ebola-neutralizing and anti-PD1

### 2

7

### clones are predisposed

- 3 Yanfang Zhang<sup>1,2,3,4,a\*</sup>, Qingxian Xu<sup>5,b\*</sup>, Huikun Zeng<sup>1,2,3,4,c\*</sup>, Minhui Wang<sup>1,6,7,d\*</sup>, Yanxia Zhang<sup>1,2,e\*</sup>,
- 4 Chunhong Lan<sup>1,3,f\*</sup>, Xiujia Yang<sup>1,2,3,4,g\*</sup>, Yan Zhu<sup>1,2,h\*</sup>, Yuan Chen<sup>3,i\*</sup>, Qilong Wang<sup>3,j</sup>, Haipei Tang<sup>3,k</sup>, Yan
- 5 Zhang<sup>2,1</sup>, Jiaqi Wu<sup>2,m</sup>, Chengrui Wang<sup>2,n</sup>, Wenxi Xie<sup>1,2,o</sup>, Cuiyu Ma<sup>1,2,p</sup>, Junjie Guan<sup>1,2,q</sup>, Shixin Guo<sup>8,r</sup>, Sen
- 6 Chen<sup>2,s</sup>, Changqing Chang<sup>9,t</sup>, Wei Yang<sup>10,u</sup>, Lai Wei<sup>8,v</sup>, Jian Ren<sup>5,w†</sup>, Xueqing Yu<sup>11†</sup>, and Zhenhai Zhang<sup>1,2,3,4,x†</sup>

1

8	<sup>1</sup> State Key	Laboratory	of Organ	Failure	Research.	National	Clinical	Research	Center fo	r Kidnev	/ Disease,
---	------------------------	------------	----------	---------	-----------	----------	----------	----------	-----------	----------	------------

- 9 Division of Nephrology, Nanfang Hospital, Southern Medical University, Guangzhou, 510515, China
- 10 <sup>2</sup>Department of Bioinformatics, School of Basic Medical Sciences, Southern Medical University, Guangzhou
- 11 510515, China
- 12 <sup>3</sup>Center for Precision Medicine, Guangdong Provincial People's Hospital, Guangdong Academy of Medical
- 13 Sciences, Guangzhou 510080, China
- <sup>14</sup> <sup>4</sup>Key Laboratory of Mental Health of the Ministry of Education, Guangdong-Hong Kong-Macao Greater Bay
- 15 Area Center for Brain Science and Brain-Inspired Intelligence, Southern Medical University, Guangzhou
- 16 510515, China
- 17 <sup>5</sup>State Key Laboratory of Oncology in South China, Cancer Center, Collaborative Innovation Center for
- 18 Cancer Medicine, School of Life Sciences, Sun Yat-sen University, Guangzhou 510060, China
- 19 <sup>6</sup>Hainan Affiliated Hospital of Hainan Medical College, Haikou, Hainkou 570311, China
- 20 <sup>7</sup>Department of Nephrology, Hainan General Hospital, Haikou 570311, China
- 21 <sup>8</sup>State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-sen University, Guangzhou
- 22 510060, China
- <sup>23</sup> <sup>9</sup>Integrate Microbiology Research Center, South China Agricultural University, Guangzhou, 510642, China
- <sup>24</sup> <sup>10</sup>Department of Pathology, School of Basic Medical Sciences, Southern Medical University, Guangzhou,
- 25 510515, China
- <sup>11</sup>Division of Nephrology, Guangdong Provincial People's Hospital, Guangdong Academy of Medical Sciences,
   Guangzhou 510080, China

28

- <sup>29</sup> \*These authors contributed equally to this work.
- 30 <sup>†</sup>To whom correspondence should be addressed:

31 Zhenhai Zhang, zhenhaismu@163.com; zhangzhenhai@gdph.org.cn

- 32 Xueqing Yu, <u>yuxueqing@gdph.org.cn</u>
- 33 Jian Ren, renjian@sysucc.org.cn
- 34 **ORCID**:

35	<sup>a</sup> 0000-0001-9309-7347	<sup>b</sup> 0000-0002-1530-5531	°0000-0001-9495-5649
36	<sup>d</sup> 0000-0001-8121-7786	°0000-0003-3623-5365	f0000-0001-5030-8247
37	g0000-0003-4036-4995	<sup>h</sup> 0000-0003-1105-6491	<sup>i</sup> 0000-0001-9043-5240
38	<sup>j</sup> 0000-0002-2248-0266	<sup>k</sup> 0000-0002-5533-7263	<sup>1</sup> 0000-0002-3681-9937
39	<sup>m</sup> 0000-0003-2204-3557	<sup>n</sup> 0000-0003-1487-0595	°0000-0001-6759-7639
40	p0000-0001-7445-6332	90000-0002-9008-9242	r0000-0001-8393-9352
41	<sup>s</sup> 0000-0002-6720-8215	<sup>t</sup> 0000-0002-5301-2932	<sup>u</sup> 0000-0001-9438-7215
42	v0000-0002-3300-8506	w0000-0002-4161-1292	×0000-0002-4310-0525

## 43 Abstract

44 Antibody repertoire refers to the totality of the superbly diversified antibodies within an individual to cope 45 with the vast array of possible pathogens. Despite this extreme diversity, antibodies of the same clonotype, 46 namely public clones, have been discovered among individuals. Although some public clones could be 47 explained by antibody convergence, public clones in naïve repertoire or virus-neutralizing clones from not 48 infected people were also discovered. All these findings indicated that public clones might not occur by 49 random and they might exert essential functions. However, the frequencies and functions of public clones in a 50 population have never been studied. Here, we integrated 2,449 Rep-seq datasets from 767 donors and discovered 5.07 million public clones  $- \sim 10\%$  of the repertoire are public in population. We found 38 51 52 therapeutic clones out of 3,390 annotated public clones including anti-PD1 clones in healthy people. Moreover, 53 we also revealed clones neutralizing SARS-CoV-2, Ebola, and HIV-1 viruses in healthy individuals. Our result 54 demonstrated that these clones are predisposed in the human antibody repertoire and may exert critical 55 functions during particular immunological stimuli and consequently benefit the donors. We also implemented RAPID – a Rep-seq Analysis Platform with Integrated Databases, which may serve as a useful tool for others 56 in the field. 57

58 Keywords: antibody repertoire, public clone, neutralizing antibody, therapeutic antibody, analysis platform

59

## 60 Background

Antibody is a critical immunoglobulin complex consisting of two identical heavy and two identical light chains. Each chain is encoded by selectively recombining one of the various germline gene fragments, namely variable (V), diversity (D, for heavy chain only), and joint (J) genes. The sequence between V gene end and J gene start is called complementarity determining region 3 (CDR3) which is extremely diverse because of the random

nucleotide insertion and deletion in the junctions and by large defines the binding specificity of an antibody<sup>1</sup>.

66 This binding specificity makes antibodies favorable for therapeutic purposes.

With tremendous efforts and various techniques, many monoclonal antibodies (mAbs) targeting distinct viruses and proteins were discovered in the past decades<sup>2,3</sup>. However, the primary barrier to studying antibodies is their immense diversity. The total number of antibodies in an adult, termed antibody repertoire, is estimated to be around  $10^{12}$  – a number far out of reach for these traditional methods<sup>4</sup>.

Fortunately, antibody repertoire sequencing (Rep-seq) was invented to acquire millions of antibody variable regions in DNA or RNA form in a single experiment, a great advance thanks to the advent of high-throughput sequencing (HTS) technology. With the aid of this technique, our understanding of the humoral immunity was markedly advanced and many valuable mAbs were identified. For instance, we and others have used Rep-seq method to discover HIV-1 broad neutralizing antibodies<sup>5-7</sup>. It also helped researchers in identifying neutralizing antibodies in the recent SARS-CoV-2 outbreak<sup>8-10</sup>. Thus far, Rep-seq has been proved to be productive in studying cancer immunology<sup>11</sup>, virus infection<sup>12,13</sup>, vaccination<sup>14,15</sup>, etc.

78 Besides the achievement aforementioned, this data-rich method also led to the finding of public clone -79 antibodies in different individuals but share the same or similar CDR3 which implies the same binding specificity. The fraction of public clones between two individuals is estimated to be ~0.95% to 6%<sup>16,17</sup> in 80 81 circulating repertoire, or linearly correlate with the product of total clones of the sample pair<sup>18</sup>. They were 82 found in individuals infected with the same virus and thus implicated the antibody convergence, a phenomenon 83 in which antibodies are assimilated to each other<sup>19</sup>. Later, they were also present in B1 and marginal zone B cells in the naïve state<sup>20</sup>. For instance, Soto et al. revealed public clones in cord blood<sup>16</sup>. Intriguingly, Jardine et 84 85 al. found VRC01-class HIV-1 neutralizing antibody clones in naïve B cells from healthy individuals<sup>21</sup>. Kreer et 86 al. discovered SARS-CoV-2-neutralizing clones in uninfected healthy people<sup>8</sup>. All these studies were 87 conducted on a limited number of samples, the answers to some of the key questions about public clone 88 remained unsolved. What proportion of an antibody repertoire is public at a population level? What are the 89 other public clones existing in the human repertoire? Have they undergone maturation process? What are their 90 functions? Do they influence our health during disease onset or virus infection?

91 Bearing these questions in mind, we collected 88,059 known antigen-binding or disease-associated antibodies 92 published before, 521 therapeutic antibodies recorded by the World Health Organization (WHO)<sup>22</sup>, and 2,449 high-quality Rep-seq datasets (767 donors, 306 million clones, and 7.12 billion raw reads) published by others 93 94 as well as generated in our lab. Integrative and systematic analysis revealed that there are around  $\sim 10\%$  or 95 more public clones for each individual in a population level. Three thousand three hundred and ninety of these 96 public clones can be annotated indicating they are functionally important for humoral immune response. More 97 importantly, we found public clones that binding to PD1, neutralizing SARS-CoV-2, Ebola, and HIV-1 viruses 98 in healthy individuals. These results demonstrated that public clones in the population are predisposed in the 99 repertoires of particular individuals who may later benefit from their existence upon virus infections and 100 disease onset. All datasets in this study were integrated and implemented in RAPID - Rep-seq Analysis 101 Platform with Integrated Database, a knowledge-rich platform for others to analyze and annotate their own 102 repertoire data.

103

### 104 **Result**

#### 105 **RAPID: a powerful platform for Rep-seq data analysis**

106 Currently, a substantial number of tools or web servers have been proposed to address the issues of Rep-seq 107 data analysis or characterization for repertoires<sup>23-40</sup>. However, these platforms focus on analyzing Rep-seq 108 dataset individually and ignoring exploration of discriminating repertoire features within or among groups. 109 Apart from that, antibody databases are also specialized for antigen annotation, such as bNAber which just 110 documents HIV broadly neutralizing antibodies<sup>41</sup>. Thus, our platform named RAPID which compensates for 111 shortages above was built. As shown in Fig. 1a, the data repositories comprised three different data modules, 112 namely Rep-seq data collection, therapeutic antibody collection, and known antibody collection. The Rep-seq data integrated 2,449 high-quality datasets (see Yang et al. for method<sup>18</sup>) from 767 donors either downloaded 113

114 from published data repository or generated in our lab. These datasets contain samples from different genders, 115 various tissues, immune status, and age spans, and were generated via different amplification strategies. Thus, 116 it provided a rich source of reference for analyzing and comparing antibody repertoires. There are 7.12 billion 117 reads and 306 million clones yielded from a systematic analysis pipeline using exactly the same criteria, thus making them comparable to each other<sup>18</sup>. The therapeutic monoclonal antibodies (mAbs) were downloaded 118 119 from the Thera-SAbDab database which contains 521 therapeutic mAbs of different types at various stages. 120 The 88,059 known antibodies were downloaded from multiple data repositories and carefully annotated via 121 natural language processing method as well as manual check (Supp. Fig. 1 and Materials and Methods). These 122 annotations included antibody sequences of different chain types as well as the antibodies that binding to and 123 neutralizing virus, associating to particular diseases, etc. The raw sequences, CDR3s, descriptions, and sample 124 metadata information were systematically extracted and stored in a rational database as well as FASTA format 125 files when necessary. A user-friendly interface for searching particular terms, antibodies, and CDR3s (Supp. 126 Fig 2) was also provided (https://rapid.zzhlab.org/).

127 The data analysis module allows users to streamline their own data through a versatile pipeline. Apart from the 128 general low- and high-level analyses, the RAPID also provides some helpful features as described below (Fig. 129 1b). 1) Customizing germline reference. 2) Customizing reference datasets; the users can freely select one or 130 more datasets in the platform as reference for cross-comparisons purpose. 3) Automatic antibody annotation; 131 the CDR3s from the input dataset will be automatically compared to the CDR3s in the data repository on 132 RAPID and annotated where applicable. 4) Downloadable figures and analysis result. Thus, any researcher can 133 upload their datasets and cross compare them to 2,449 datasets with 306 million clones, all therapeutic 134 antibodies, and 88,059 known antibodies and retrieve the relevant information. With the thorough antibody 135 collections and a versatile analysis platform, we believe the RAPID will be helpful for the large cadre of 136 scientists who demand analyzing antibody repertoire data as we demonstrated in this study.

137

138

139



b



140

141 Figure 1 Data composition and automatic online Rep-seq dataset analysis pipeline of the RAPID. (a) 142 Composition and detailed information of the three types of antibody datasets: Rep-seq datasets, therapeutic 143 antibodies and known antibodies. The pie charts in the lower panel showed detailed composition of each type 144 of antibodies. (b) The analysis platform for Rep-seq dataset.

145

#### 146 **Public clones are prevalent in the population**

147 With this unprecedented dataset, we started in-depth inspection of the public clones. Here, we defined public

148 clones as antibodies with the same CDR3 amino acid sequence that present in more than one donor. Even with

149 this stringent criterion, we discovered 5,077,372 public clones. Typically, the public clones represent ~1.23 to 150 100 percent of each individual repertoire with a peak value of 10.46 percent (Fig. 2a). Moreover, 65 public 151 clones occur in more than 100 individuals with one clone shared by 196 individuals (25.55% of the total 152 donors) (Supp. Fig. 3 a and b). Thus, population level study helped us find more public clones and highly 153 frequent ones. As 96.86% of the public clones are from PBMCs (Supp. Fig. 3c), we compared their SHMs and 154 clone fractions of the clones in different groups. The SHMs for naïve, memory, and plasma groups were 155 comparable between public clones and total clones. The public clones of PBMCs and unknown samples 156 displayed mediocre SHMs between naïve and non-naïve clones (Fig. 2b, upper panel). For clone fractions, the 157 public clones from PBMCs and unknown samples were lower than the other counterparts (Fig. 2b, lower panel) 158 indicating they are inactivated. On the other hand, about half of the public clones were from IgM isotype (Fig. 159 2c, lower panel). Therefore, it's reasonable to speculate that majority of these public clones were acquired 160 from naïve and lowly-mutated memory B cells.

161 We also observed that different V and J gene combinations can yield the same CDR3s. As expected, the 162 diversity of V gene usage for public clones increased when clones are shared by more donors (Supp. Fig. 4). 163 However, when normalized to the maximum theoretical diversity with particular number of V genes (see 164 Materials and Methods), this diversity slightly decreased with the increment of sharing donors indicating 165 recombination preference of V genes (Fig. 2c, upper panel). Careful examinations of the V and J genes that 166 formed the same CDR3s showed that same J gene was always preferred while V gene was more replaceable 167 among individuals (Fig. 2d and Supp. Fig. 5). Nevertheless, the substitution rates of V genes were not 168 completely influenced by their sequence similarity (Fig. 2d). This result suggested that J genes might affect the 169 CDR3s more than V genes do.

Taken advantage of the rich antibody information integrated in this study, we tried to annotate these public clones. Totally, 3,390 public clones have been annotated by three antibody databases including known antibody, Thera-SAbDab, and Coronavirus-neutralizing antibody incorporated with 459 mAbs from CoV-AbDab<sup>42</sup>, 28 mAbs from Kreer *et al.*<sup>8</sup>, and 19 mAbs from Liu *et al.*<sup>10</sup> (Fig. 2e). We found that 3,349 out 3,390



175

176 Figure 2 The characteristics of public clones. (a) The sample density distribution with regard to the 177 fractions of public clones. X-axis represents the percent of public clones in a sample. Y-axis shows the sample 178 density. (b) Comparisons of the distribution of somatic hypermutation rates (SHMs, upper panel) and clone 179 fractions (lower panel) among different groups of clones. (c) Upper panel: The normalized Shannon index of V 180 gene usage (Y-axis) within each CDR3 aa group sorted by the respective numbers of shared donors (X-axis). Lower panel: The stacked bar plot showed the composition of public clones for different antibody isotypes. (d) 181 The substitution frequencies (lower left panel) and percent identities of V genes (upper right panel). (e) The 182 183 overlap of 3,390 annotated public clones with known antibodies, thera-SAbDab, and Coronavirus-neutralizing 184 antibodies.

185

associating with diseases. Surprisingly, we found 38 and 27 CDR3s overlapping with Thera-SAbDab and Coronavirus-neutralizing antibodies. Among 27 Coronavirus-binding clones, 8 of them target SARS-CoV-2 (CARGDSSGYYYYFDYW binds both SARS-CoV-1 and SARS-CoV-2). As all these Rep-seq datasets were generated before the outbreak of the COVID-19, these SARS-CoV-2-neutralizing clones and therapeutic clones were deposited in the antibody repertoire. These provide evidence that those therapeutic and antigenspecific antibodies exist widely in populations and can be a powerful source for mAb discovery with clinical purposes.

### 193 Therapeutic mAb clones are prevalent in healthy people's repertoire

194 We went on to look into the details of the 3,390 annotated public clones. Among them, 3,354 (98,94%) were 195 found in at least one healthy sample (Fig. 3a, upper panel). The therapeutic antibodies found in this data 196 collection compromise 41 therapeutic mAbs with 38 unique CDR3s. Of these, six are under phase III clinical 197 trial, two are under preregistration and nine are approved for clinical usage (Fig. 3a, lower panel). Interestingly, 198 only 14 (34.15%) mAbs are fully human while the others include 14 (34.15%) humanized, 12 (29.27%) from 199 mouse, and one chimeric. This indicated that therapeutic antibodies generated by mouse model can be 200 generated by our own immune system. Thus, public clone might be a better source for mAb discovery for 201 clinical usage.

202 Also, many of the therapeutic antibody clones found in healthy people which are used for treatment of diseases 203 with top causes of death in the world (Table 1) prevailed in the population (Fig. 3, a and b). For instance, 204 Evolocumab targeting PCSK9 is used to treat Coronary disorders, Stroke, and Hypercholesterolaemia. Stroke alone caused 5.78 million deaths worldwide in 2016<sup>43</sup> and this clone was found in 108 (14.1% of the total of 205 206 767) donors' repertoire. The CDR3 of anti-PD1 (Camrelizumab), the treatment to various cancers, was found 207 in 14 donors' repertoire. Ramucirumab targeting KDR and Enfortumab targeting PVRL4 were also found in 23 208 and 49 donors, respectively. According to the percent identities to therapeutic antibodies, most of the 209 antibodies from the same clonotype separated into at least two groups (Fig. 3b and Supp. Fig. 6). Detailed 210 inspection revealed that multiple V genes involved in the recombination, again supported the diversity of V 211 genes within clones (Fig. 2d). These antibodies might serve as therapeutic alternatives for the same disease.



#### 212

213 Figure 3 The characteristics of annotated public clones. (a) Overview of 3.390 annotated public clones. 214 The upper heatmap stand for the composition of samples for these annotated public clones. Samples were 215 divided into 6 groups including allergy, autoimmune, cancer, pathogen, healthy, and others. The bottom line 216 chart means the number of donors. Public clones annotated by known antibody database were shown by gray 217 line, while those annotated by other databases were shown by scatters filled in particular colors. The center top 218 pie charts show distribution of source and clinical trial for therapeutic antibodies and associated disease for known antibodies. (b) Identity of variable region sequences from FR1 to FR3 of antibodies whose CDR3aa are 219 220 same as Evolocumab, Camerelizumab, Ramucirumab, and Enfortumab. The X-axis means the germline 221 divergence and different colors of scatters mean different V genes. Numbers filled in red stand for the death 222 caused by diseases treated by such antibody and those filled in blue stand for the number of variable regions 223 identified from Rep-seq datasets. Titles for sub-figures separated by forward slash include inn id of therapeutic 224 antibody, the number of samples and donors with such CDR3aa, and target of therapeutic antibody. Dots for 225 therapeutic antibodies are larger than that of clones identified from Rep-seq datasets. (c) Multiple sequences 226 alignment of variable region for antibodies with the same CDR3aa as Camerelizumab from homo sapiens, mus 227 musculus, and rattus norvegicus. The amino acid sequence of Camerelizumab is listed in the top and regions 228 (FR1-FR4) for variable region are marked by boxes above. V gene used in each sequences are labeled in left 229 boxes.

### Table 1 Detailed information of therapeutic antibodies whose CDR3aa were same as those identified from Rep-seq datasets

INN	CDR3aa	# Donor	# Sample	Source	Target	Disease	Death_rate	Total_death_rat	e # Total_death (k)
Evolocumab	CARGYGMDVW	108	345	fully human	PCSK9	Coronary disorders,Hypercholesterolaemia,Hyperlipidaemia,Hyperlipoproteinaemia typ IIa.Myocardial infarction,Stroke	-,-,-,-,10.16 •	10.16	5778.38
Azintuxizumab	CARDRGYYFDYW	62	248	chimeric	SLAMF7	-	-	-	-
Crenezumab	CASGDYW	58	181	humanized	APP	Alzheimer's disease	3.5	3.5	1990.58
Solanezumab	CASGDYW	58	181	humanized	APP	Alzheimer's disease	3.5	3.5	1990.58
Landogrozumab	CARLPDYW	56	189	humanized	MSTN	Cachexia, Muscular atrophy, Pancreatic cancer	-,-,0.65	0.65	369.68
Enfortumab	CARAYYYGMDVW	49	144	fully human	PVRL4	Urogenital cancer	-	-	-
Dusigitumab	CARDPYYYYYGMDVW	27	89	fully human	IGF1&IGF2	•	-	-	-
Eptinezumab	CARGDIW	26	67	humanized	CALCA&CALCB	Migraine	0	0.00	0.00
Ramucirumab	CARVTDAFDIW	23	78	fully human	KDR	Biliary cancer, Carcinoid tumour, Colorectal cancer, Gastric cancer, Live cancer, Non-small cell lung cancer, Oesophageal cancer, Pancreatic cancer, Soli tumours, Urogenital cancer	ei-,-,- c,1.34,1.46,- .0.75.0.65	4.2	2388.70
Daclizumab	CARGGGVFDYW	22	70	humanized	IL2RA	Multiple sclerosis, Renal transplant rejection	0.04,-	0.04	22.75
Bimagrumab	CARGGWFDYW	20	47	fully human	ACVR2B	Cachexia, Muscular atrophy, Type 2 diabetes mellitus	-,-,2.81	2.81	1598.15
Glembatumumab	CARGYNWNYFDYW	19	67	fully human	GPNMB	-	-	-	•
Zanolimumab	CARVINWFDPW	18	67	fully human	CD4		-	-	-
Valanafusp	CAREWAYW	15	54	mouse	INSR	Mucopolysaccharidosis I	-	-	-
Camrelizumab	CARQLYYFDYW	14	30	humanized	PDCD1	Biliary cancer,Breast cancer,Cervical cancer,Cholangiocarcinoma,Colorect: cancer,Diffuse large B cell lymphoma,Endometrial cancer,Fallopian tub cancer,Gastric cancer,Hodgkin's disease,Liver cancer,Malignan melanoma,Nasopharyngeal cancer,Non small cell lung cancer,Oesophagea cancer,Osteosarcoma,Ovarian cancer,Pancreatic cancer,Peritoneal cancer,Ren cell carcinoma Solid tumours Lurgential cancer	al-,1.03,-,-,-,- ic,-,1.34,- nt,1.46,-,-,- al,0.75,- a',0.29,0.65,- 0.28	5.8	3298.68
Ascrinvacumab	CARESVAGFDYW	14	37	fully human	ACVRL1	Age-related macular degeneration, Cancer, Liver Cancer	-,-,1.46	1.46	830.36
Cobolimab	CASMDYW	12	23	mouse	HAVCR2	Non-small cell lung cancer,	-,-	-	-
Reslizumab	CAREYYGYFDYW	10	17	humanized	IL5	Asthma,Churg-Strauss syndrome,Sinusitis	0.73,-,-	0.73	415.18
Tesidolumab	CARDTPYFDYW	10	15	fully human	C5	Paroxysmal nocturnal haemoglobinuria	-	-	-
Prasinezumab	CARGGAGIDYW	9	39	humanized	SNCA	Parkinson's disease	0.38	0.38	216.12
Gatralimab	CTPIDYW	9	15	mouse	CD52	-	-	-	-
Utomilumab	CARGYGIFDYW	8	18	fully human	TNFRSF9	B-cell lymphoma,Breast cancer,Colorectal cancer,Diffuse large B ce lymphoma,Follicular lymphoma,Oropharyngeal cancer,Ovarian cancer,Soli tumours	ll-,1.03,-,-,-,- d,0.29,-	1.32	750.73
Bemarituzumab	CARGDFAYW	7	18	humanized	FGFR2	Gastric cancer, Oesophageal cancer, Solid tumours	1.34,0.75,-	2.09	1188.66
Bersanlimab	CARYSGWYFDYW	6	14	mouse	ICAM1	-	-	-	-

Cont. Table 1									
INN	CDR3aa	# Donor	# Sample	Source	Target	Disease	Death_rate	Total_death_rat	te # Total_death (k)
Iratumumab	CASLTAYW	5	6	fully human	TNFRSF8	-	-	-	-
Mosunetuzumab	CARDSYSNYYFDYW	5	5	humanized	CD3E, MS4A1	Chronic lymphocytic leukaemia,Diffuse large B cell lymphoma,Non-Hodgkin	l's-,-,-	-	-
Emibetuzumab	CARANWLDYW	4	4	humanized	MET	Cancer,Gastric cancer,Non-small cell lung cancer	-,1.34,-	1.34	762.11
Otilimab	CARGFGTDFW	4	7	mouse	CSF2	Rheumatoid arthritis	0.1	0.10	56.87
Vorsetuzumab	CARDYGDYGMDYW	4	6	humanized	CD70	•	-	-	-
Gatipotuzumab	CTRHYYFDYW	3	9	humanized	MUC1	Ovarian cancer, Solid tumours	0.29,-	0.29	164.93
Pankomab	CTRHYYFDYW	3	9	mouse	MUC1	Ovarian cancer,Solid tumours	0.29,-	0.29	164.93
Brodalumab	CARRQLYFDYW	3	5	fully human	IL17RA	Erythrodermic psoriasis,Plaque psoriasis,Psoriasis,Psoriatic arthritis,Pustul psoriasis,Spondylarthritis,Systemic scleroderma	a1-,-,-,-,-,-,-	-	-
Clervonafusp	CARRGLLLDYW	3	3	mouse	SLC29A2	Glycogen storage disease type II	-	-	-
Dectrekumab	CARLWFGDLDAFDIW	3	3	fully human	IL13	•	-	-	-
Ivuxolimab	CARESGWYLFDYW	2	3	mouse	TNFRSF4	Acute myeloid leukaemia,Breast cancer,Follicular lymphoma,Renal ce	ell-,1.03,- 0.28	1.31	745.05
Cinpanemab	CTSAHW	2	2	mouse	SNCA	Parkinson's disease	0.38	0.38	216.12
Atoltivimab	CARNWNLFDYW	2	2	mouse	Zaire Ebolavirus G	P -	-	-	-
Dinutuximab	CVSGMEYW	2	2	mouse	Ganglioside GD2	Neuroblastoma,Small cell lung cancer	-,-	-	-
Lorukafusp	CVSGMEYW	2	2	mouse	Ganglioside GD2	Malignant melanoma, Neuroblastoma	-,-	-	-
Lupartumab	CAREGLWAFDYW	2	4	fully human	LYPD3	Solid tumours	-	-	-
Tildrakizumab	CARGGGGFAYW	2	3	humanized	IL23A	Ankylosing spondylitis,Intervertebral disc degeneration,Non-radiographic axi spondyloarthritis,Plaque psoriasis,Psoriatic arthritis	al-,-,-,-	-	-

233 Note: Cells filled in gary mean that there are no applied disease for these therapeutic antibodies.

Moreover, for therapeutic clones, there were always Rep-seq captured antibodies showing higher SHM rates (Fig. 3b), possibly indicated the engagement of maturation process under the selective pressure of their antigens.

237 Interestingly, we also found antibodies with the same CDR3s as anti-PD1 antibody in the repertoires of mouse 238 and rat. The variations in the other regions of variable sequences (Fig. 3c) indicated that they were purposely 239 selected and retained in the repertoire. Moreover, the differences in their structure might be the result of 240 structural variations of PDCD1 proteins in different species (Supp. Fig. 7). Finding anti-PD1 antibody clones 241 in healthy individuals and other species was exciting. On the contrary, no anti-PD1 clone was found in any of 242 the 56 samples of various cancers -2 breast cancer, 49 colorectal cancer, and 5 liver cancer samples. Why 243 cancer patients do not possess anti-PD1 antibody clones is an intriguing question. Understanding the 244 mechanism behind this fact may lead to a better understanding of cancer immunology and more effective 245 immunological therapy. Besides, further studies illustrating the causes and consequences of the anti-PD1 246 clones in particular individuals would also benefit the field.

### 247 SARS-CoV-2-, Ebola-, and HIV-1-neutralizing antibody clones are predisposed

On a par with the results in finding therapeutic clones, multiple clones neutralizing SARS-CoV-2, Ebola, and 248 249 HIV-1 were also uncovered from virus-naïve individuals. For SARS-CoV-2-neutralizing clones MT658807, 250 MT658819, and 1-20, there were 20, 506, and 222 heavy chain variable regions extracted from repertoires of 5, 251 17, and 6 healthy donors, respectively. Twenty-five variable regions from 2 healthy donors sharing the same 252 CDR3 sequence with HIV-1-neutralizing class VRC01 were also extracted. In addition, 2,663 variable regions 253 from 4 donors injected with influenza vaccine contained the CDR3 of Ebola-neutralizing antibody of 254 MK90182344. Although neutralizing clones to HIV-1 and SARS-CoV-2 were reported to exist in the repertoire 255 of the naïve B cell previously<sup>8,21</sup>, this is the first time to identify these neutralizing clones in multiple people. 256 Thus, we concluded that they are predisposed in a population.

We then set off to explore the maturation pathways of these neutralizing clones by analyzing the phylogenetic trees of each clone built via DNAMLK (Fig. 4a-d). The Ebola-neutralizing clones exhibited high maturation rates with IgG. Interestingly, the maturation rates of three SARS-CoV-2 clones demonstrated various level of

SHMs. While the overall SHM rate for MT658807 clone is lower than 2.5%, some of the antibodies in MT658819 and 1-20 clones displayed more than 5% mutations. Previous studies reported the general lower SHM for SARS-CoV-2-neutralizing antibodies but still some clones with more mutations were identified and verified<sup>8-10</sup>. Therefore, different clones might subject to different selective pressure and consequently manifest various SHM rates. Although the diversity of V gene usage in neutralizing clones by large defined the topology of the phylogenetic trees, oftentimes antibodies from different individuals aggregating in the same branch was observed. This indicated antibody convergence in the same maturation pathway.

Apart from the MK901823 clone, which was from a sample after influenza vaccine trivalent, inactivated seasonal influenza (TIV), all the donors of the HIV-1 and SARS-CoV-2 clones are virus-free healthy individuals. Furthermore, we found the same antibodies of MT658807 in donor 1776. This again confirmed the predisposition of this neutralizing clone. However, what triggers their maturation would be an important question to answer for future studies.

To further explore their possibility of virus binding, we compared the structures of these Rep-seq retrieved antibodies to their corresponding verified neutralizing antibodies. As shown in Fig. 4a-d, we found antibodies were very similar to MT658819<sup>8</sup> (RMSD: 0.165) and 1-20<sup>10</sup> (RMSD: < 0.001) that neutralizing SARS-CoV-2,

275 MK901823<sup>44</sup> (RMSD: 0.802) that neutralizing Ebola, and KU760937<sup>21</sup> (RMSD: 0.459) that neutralizing HIV.

To validate this similarity, we performed pair-wise structure comparison among antibodies the neutralizing these three viruses. As shown in Fig. 4e, the RMSD scores of clone targeting the same antigen were much lower than those targeting different antigens. Thus, the high similarity of Rep-seq retrieved antibodies to neutralizing antibodies are reliable.

### 280 SARS-CoV-1-neutralizing and therapeutic antibody clones exist in animals

Inspired by the existence of anti-PD1 clones in mouse and rat, we scrutinized the Rep-seq datasets with four different species, namely Macaca fascicularis, Macaca mulatta, Mus musculus, and Rattus norvegicus. We found 4 SARS-CoV-1-neutralizing and 18 therapeutic clones in at least one species. Taken together, we believe these clones are not randomly generated but purposely selected and disposed in vertebrates' repertoire.



285

Figure 4 Maturation pathway and structure of antibodies potentially targeting SARS-CoV-2, Ebola, 286 287 and HIV. Variable region sequences with the same CDR3aa as MT658819 (a), 1-20 (b), MK901823 (c), and 288 KU760937 (d) were extracted and compared with those validated neutralizing antibodies and their germline 289 reference. The germline reference was chosen as root of phylogenetic tree and the validated antibodies were 290 marked by arrow. The cluster map contains four layers including similarity of sequences (the sequences 291 extracted from the same donor were marked with the same color), V gene family, isotype, and somatic 292 hypermutation rate from inner to outer. Under each phylogenetic tree, similarity of structures for validated 293 antibody and variable region identified from Rep-seq datasets were shown. (e) Pair-wise structure comparison 294 of antibodies targeting SARS-CoV-2, HIV, and Ebola.

# 295 **Discussions**

Public clones are a specific fraction of antibodies among individuals that we know little about. By integrating the largest antibody data to date, population level analyses discovered millions of public clones which represent ~10% or higher fraction of each individual's repertoire. However, compared to the superb diversity of the antibody repertoire, the current dataset might still be smaller than demand. We believe that when more datasets will be integrated, more public clones would be revealed. This is understandable since although the somatic recombination may generate numerous antibodies, majority of them are eliminated during the negative selection process in the bone marrow. Consequently, the once private repertoire might be public<sup>45</sup>.

303 How often can we find these public clones with critical functions in an individual? Are they predisposed in 304 everyone's repertoire? The current data seems to support that only some people possess them. However, we 305 found that sequencing depth is critical for public clone identification as many more public clones were 306 observed in datasets with very high depth. Currently, only a few hundred thousand to a few million reads were 307 captured in general. Compared to the theoretical number of B cells in the sample and the depth needed to 308 identify a clone confidently, much more sequencing reads are demanded. As most of the therapeutic mAbs 309 target proteins of conserved genes such as PDCD1, another helpful practice in finding functional public clones 310 might be comparing antibody repertoires between human and other vertebrates.

The finding of clones that can bind to PDCD1 or neutralize SARS-CoV-2, Ebola, and HIV-1 viruses demonstrated that public clones might be important for the donor's health. Then discovering the functionalities of the vast majority of other public clones would be critical for a deep understanding of the humoral immune system. The major challenge in this regard is the lack of the light chain pair. The techniques of paired heavy and light chain sequencing invented in Georgiou lab<sup>46</sup> and the single cell repertoire sequencing<sup>47</sup> showed great potential in solving this problem.

We'll update RAPID along with the accumulation of Rep-seq datasets generated by others and our lab. Webelieve more public clones will be identified and their functions will be illustrated along this path.

## 319 Methods

#### 320 Rep-seq datasets enrollment

- 321 Method to enroll published and in-house Rep-seq datasets were described in Yang *et al*<sup>18</sup>, please refer to it for
- 322 detailed information. The re-analysis pipeline of these Rep-seq datasets was also included in that paper.

#### 323 **Resources of known antibody**

Five open access antibody databases, named abYsis (http://abysis.org/)<sup>48</sup>, bNAber (http://bnaber.org/)<sup>41</sup>, 324 325 EMBLIG (http://acrmwww.biochem.ucl.ac.uk/abs/abybank/emblig/), HIV Molecular Immunology Database (HIV-DB: https://www.hiv.lanl.gov/content/immunology/neutralizing ab resources.html)<sup>49</sup> and IMGT/LIGM-326 DB (http://www.imgt.org/ligmdb/)<sup>50</sup> were enrolled. In addition, another two nucleotide sequence databases, 327 328 including European Nucleotide Archive (ENA) of EMBL-EBI 329 (https://www.ebi.ac.uk/ena/data/view/Taxon:9606&result=coding\_release)<sup>51</sup> and National Center for 330 Biotechnology Information (NCBI) Nucleotide database (https://www.ncbi.nlm.nih.gov/nucleotide/), were also 331 incorporated. In a word, 7 databases were finally included. The search strategy and download date for them 332 were listed in Supplementary Table 1.

#### 333 Construction of known antibody database

Although the species was restricted for sequences downloading, some sequences from other species, like Mus 334 335 musculus, were also included. Thus, we firstly discarded non-human sequences according to descriptions. 336 After that, sequences were aligned to V, D, and J germline reference (downloaded from IMGT: 337 http://www.imgt.org/ and listed in Supplementary Table 2) by IgBLAST<sup>24</sup> (version 1.8.0), as its great 338 performance for error-corrected reads<sup>23</sup>. Based on results of IgBLAST, sequences which meet criteria were 339 reserved, including in-frame, productive, with V, J, and CDR3, without either stop codon or out-of-frame in 340 variable region, and without ambiguous base (N) in CDR3. Sequences with the same nucleotide sequences of variable region within the same database were de-duplicated. To remove non-antibodies from NCBI and ENA, 341 342 we aligned these sequences to NCBI Nt database (downloaded from NCBI:

343 ftp://ftp.ncbi.nlm.nih.gov/blast/db/FASTA/nt.gz) and discarded the sequences whose descriptions contain no 344 keywords we defined (Listed in Supplementary Table 3). These keywords were selected from descriptions of 345 the antibody sequences are stored in the database of abYsis, bNAber, HIV-DB, EMBLIG, and IMGT/LIGM-DB. Furthermore, antibodies from 7 databases were pooled together and de-duplicated according to the 346 347 nucleotide sequence of variable region. In the end, disease information for antibodies from EMBLIG, ENA, 348 IMGT/LIGM-DB, and NCBI was annotated by TaggerOne<sup>52</sup> based on description, title, and abstract of 349 sequences. The sequences from abYsis were annotated as "NA", as no annotation information can be 350 downloaded. The sequences from HIV-DB and bNAber were annotated as HIV infections.

#### 351 Implementation of RAPID

The web interface is implemented by Hyper Text Markup Language (HTML), Cascading Style Sheets (CSS), and JavaScript (JS). It is a single page application based on the JS framework React.js, while using the React component library Ant Design to unify the design style. The back end of the website uses Nginx as the HTTP and reverse proxy server, develops business logic based on Node.js, uses MySQL to manage data, and uses RabbitMQ to process the analysis task queues. Furthermore, the real-time notification of task progress depends on the WebSocket technology.

#### 358 Extraction of variable region identified from Rep-seq dataset

Firstly, if regions from FR1 to FR4 were reported by MiXCR, we would simply join them together as variable region. For sequences whose FR1 to FR4 regions were not completely reported by MiXCR, we extracted them using our algorithm: I) Reads which can not be merged by MiXCR were discarded; II) The beginning of variable region was acquired by pairwise alignment between germline reference of V genes and the column named "targetSequence" reported by MiXCR(The function *pairwise2.align.localms* from Python *Bio module* was used with parameters 2, -3, -5, and -2); III) If the column named "refPoints" in MiXCR recorded the region of FR4, we would use it instead of aligning "targetSequence" to J gene to find the end of FR4..

366

367

#### 368 Calculation of gene usage diversity

The Shannon index was used to show the diversity of gene usage of public clones. However, as the number of V genes influences the diversity largely, we used the maximum of diversity with particular number of V genes to normalized the diversity. The function to calculate the normalized diversity is shown below. When different donors use totally different V genes, the normalized diversity equals one.

373

Normalized V gene diversity = 
$$\frac{\sum_{i=1}^{N} (P_i * \ln P_i)}{\ln N}$$

374  $P_i$  means the frequency of specific V gene, *i* means the order of V gene, and *N* means the total number of V 375 gene.

#### 376 Calculation of sequences identity

Both nucleotide and amino acid sequences were aligned by Clustal W 2.1<sup>53</sup> from the Python module named
Bio. Gaps at the beginning and ending of aligned sequences were removed and the percentage of matched
bases was defined as identity.

### 380 Multiple sequence alignment for anti PD-1 antibodies

Variable region sequences with the same CDR3aa of Camrelizumab were extracted from each sample and grouped according to the VJ recombination and CDR3nt. Amino acid sequences of groups with most reads in each individual were used for multiple sequence alignment by Clustal W 2.1 with default parameters and visualized by BioEdit.

### 385 Construction of phylogenetic tree

Each phylogenetic tree was generated by the nucleotide sequences of variable regions for antibodies sharing the same CDR3 sequence with MT658807, MT658819, 1-20, MK901823, and KU760937. In addition, the germline V allele of validated neutralizing antibody which was set as the root and validated antibody were also

enrolled. Alignments were performed using Clustal W 2.1, and the maximum parsimony trees fitted using
 DNAMLK by PHYLIP 3.698<sup>54</sup>. Lastly, these phylogenetic trees were displayed and annotated by iTOL<sup>55</sup>.

#### 391 Comparison of antibody structure

392 As some Rep-seq datasets were amplified by Multiplex PCR, variable regions for these sequences were not

393 complete. Thus, sequences lost several bases at the beginning of the FR1 due to the design of primer set were

394 padded by germline sequences from IMGT. Sequences for validated antibodies were downloaded from NCBI.

395 Variable regions without out-frame were used to predict their structures by Repertoire Builder<sup>56</sup>. Then PyMOL

396 was used to calculate RMSD to compare the similarity of antibody structures.

397

# 398 **Reference**

- Xu, J. L. & Davis, M. M., Diversity in the CDR3 Region of VH Is Sufficient for Most Antibody
  Specificities. *Immunity (Cambridge, Mass.)* 13 37 (2000).
- Lai, J. Y. & Lim, T. S., Infectious disease antibodies for biomedical applications: A mini review of
   immune antibody phage library repertoire. *Int J Biol Macromol* 163 640 (2020).
- Frenzel, A., Schirrmann, T. & Hust, M., Phage display-derived human antibodies in clinical
   development and therapy. *Mabs-Austin* 8 1177 (2016).
- <sup>4</sup> Bruce Alberts, A. J. J. L., *The Generation of Antibody Diversity*. (Garland Science, New York, 2002).
- Wu, X. *et al.*, Maturation and Diversity of the VRC01-Antibody Lineage over 15 Years of Chronic
  HIV-1 Infection. *Cell* 161 470 (2015).
- <sup>6</sup> Setliff, I. *et al.*, High-Throughput Mapping of B Cell Receptor Sequences to Antigen Specificity. *Cell* **179** 1636 (2019).
- 410 <sup>7</sup> Setliff, I. *et al.*, Multi-Donor Longitudinal Antibody Repertoire Sequencing Reveals the Existence of
  411 Public Antibody Clonotypes in HIV-1 Infection. *Cell Host Microbe* 23 845 (2018).

<sup>8</sup> Kreer, C. *et al.*, Longitudinal Isolation of Potent Near-Germline SARS-CoV-2-Neutralizing
Antibodies from COVID-19 Patients. *Cell* (2020).

414 <sup>9</sup> Cao, Y. *et al.*, Potent Neutralizing Antibodies against SARS-CoV-2 Identified by High-Throughput
415 Single-Cell Sequencing of Convalescent Patients' B Cells. *Cell* 182 73 (2020).

- 416 <sup>10</sup> Liu, L. *et al.*, Potent neutralizing antibodies directed to multiple epitopes on SARS-CoV-2 spike.
- 417 *Nature* (2020).
- 418 <sup>11</sup> Zhang, W. et al., Characterization of the B Cell Receptor Repertoire in the Intestinal Mucosa and of
- 419 Tumor-Infiltrating Lymphocytes in Colorectal Adenoma and Carcinoma. J Immnuol 198 3719 (2017).

420 <sup>12</sup> Roskin, K. M. *et al.*, Aberrant B cell repertoire selection associated with HIV neutralizing antibody

- 421 breadth. *Nat Immunol* **21** 199 (2020).
- 422 <sup>13</sup> McCarthy, K. R. *et al.*, Memory B Cells that Cross-React with Group 1 and Group 2 Influenza A
  423 Viruses Are Abundant in Adult Human Repertoires. *Immunity* 48 174 (2018).
- 424 <sup>14</sup> Joyce, M. G. *et al.*, Vaccine-Induced Antibodies that Neutralize Group 1 and Group 2 Influenza A
  425 Viruses. *Cell* 166 609 (2016).
- Jackson, K. J. L. *et al.*, Human Responses to Influenza Vaccination Show Seroconversion Signatures
   and Convergent Antibody Rearrangements. *Cell Host Microbe* 16 105 (2014).
- 428 <sup>16</sup> Soto, C. *et al.*, High frequency of shared clonotypes in human B cell receptor repertoires. *Nature* 566
  429 398 (2019).
- <sup>17</sup> Briney, B., Inderbitzin, A., Joyce, C. & Burton, D. R., Commonality despite exceptional diversity in
  the baseline human antibody repertoire. *Nature (London)* 566 393 (2019).
- 432 <sup>18</sup> Yang, X. *et al.*, Large-scale Analysis of 2,152 dataset reveals key features of B cell biology and the
  433 antibody repertoire. *bioRxiv* 814590 (2019).
- <sup>19</sup> Parameswaran, P. *et al.*, Convergent Antibody Signatures in Human Dengue. *Cell Host Microbe* 13
  691 (2013).
- <sup>20</sup> Bautista, D. *et al.*, Differential Expression of IgM and IgD Discriminates Two Subpopulations of
  Human Circulating IgM+IgD+CD27+ B Cells That Differ Phenotypically, Functionally, and Genetically. *Front Immunol* 11 (2020).

- Jardine, J. G. *et al.*, HIV-1 broadly neutralizing antibody precursor B cells revealed by germlinetargeting immunogen. *Science (American Association for the Advancement of Science)* **351** 1458 (2016).
  Raybould, M. I. J. *et al.*, Thera-SAbDab: the Therapeutic Structural Antibody Database. *Nucleic Acids Res* **48** D383 (2020).
  Zhang, Y. *et al.*, Tools for fundamental analysis functions of TCR repertoires: a systematic
- 444 comparison. Brief Bioinform (2019).
- 445 <sup>24</sup> Ye, J., Ma, N., Madden, T. L. & Ostell, J. M., IgBLAST: an immunoglobulin variable domain
  446 sequence analysis tool. *Nucleic Acids Res* 41 W34 (2013).
- <sup>25</sup> Zhang, W. *et al.*, IMonitor: A Robust Pipeline for TCR and BCR Repertoire Analysis. *Genetics* 201
  448 459 (2015).
- Kuchenbecker, L. *et al.*, IMSEQ a fast and error aware approach to immunogenetic sequence
  analysis. *Bioinformatics* **31** 2963 (2015).
- 451 <sup>27</sup> Yu, Y., Ceredig, R. & Seoighe, C., LymAnalyzer: a tool for comprehensive analysis of next
  452 generation sequencing data of T cell receptors and immunoglobulins. *Nucleic Acids Res* 44 e31 (2016).
- 453 <sup>28</sup> Bolotin, D. A. *et al.*, MiXCR: software for comprehensive adaptive immunity profiling. *Nat Methods*454 **12** 380 (2015).
- 455 <sup>29</sup> Gerritsen, B., Pandit, A., Andeweg, A. C. & de Boer, R. J., RTCR: a pipeline for complete and 456 accurate recovery of T cell repertoires from high throughput sequencing data. *Bioinformatics* **32** 3098 (2016).
- <sup>30</sup> Hung, S. *et al.*, TRIg: a robust alignment pipeline for non-regular T-cell receptor and immunoglobulin
   sequences. *BMC Bioinformatics* 17 (2016).
- <sup>31</sup> Bolotin, D. A. *et al.*, MiTCR: software for T-cell receptor sequencing data analysis. *Nat Methods* 10
  813 (2013).
- <sup>32</sup> Alamyar, E., Giudicelli, V., Li, S., Duroux, P., & Lefranc, M. P., IMGT/HIGHV-QUEST: THE IMGT
  WEB PORTAL FOR IMMUNOGLOBULIN (IG) OR ANTIBODY AND T CELL RECEPTOR (TR)
  ANALYSIS FROM NGS HIGH THROUGHPUT AND DEEP SEQUENCING. *Immunome Res* 8 (2012).

- 464 <sup>33</sup> Thomas, N., Heather, J., Ndifon, W., Shawe-Taylor, J. & Chain, B., Decombinator: a tool for fast,
- 465 efficient gene assignment in T-cell receptor sequences using a finite state machine. *Bioinformatics* 29 542
  466 (2013).
- 467 <sup>34</sup> Yang, X. *et al.*, TCRklass: A New K-String Based Algorithm for Human and Mouse TCR
- 468 Repertoire Characterization. *J Immnuol* **194** 446 (2015).
- 469 <sup>35</sup> Zhang, W. *et al.*, PIRD: Pan Immune Repertoire Database. *Bioinformatics* **36** 897 (2020).
- 470 <sup>36</sup> Avram, O. *et al.*, ASAP A Webserver for Immunoglobulin-Sequencing Analysis Pipeline. *Front*471 *Immunol* 9 (2018).
- 472 <sup>37</sup> Christley, S. et al., VDJServer: A Cloud-Based Analysis Portal and Data Commons for Immune
- 473 Repertoire Sequences and Rearrangements. Front Immunol 9 (2018).
- 474 <sup>38</sup> Margreitter, C. *et al.*, BRepertoire: a user-friendly web server for analysing antibody repertoire data.
- 475 *Nucleic Acids Res* **46** W264 (2018).
- <sup>39</sup> IJspeert, H. *et al.*, Antigen Receptor Galaxy: A User-Friendly, Web-Based Tool for Analysis and
  Visualization of T and B Cell Receptor Repertoire Data. *J Immnuol* **198** 4156 (2017).
- 478 <sup>40</sup> Duez, M. *et al.*, Vidjil: A Web Platform for Analysis of High-Throughput Repertoire Sequencing.
  479 *Plos One* **11** e166126 (2016).
- 480 <sup>41</sup> Eroshkin, A. M. *et al.*, bNAber: database of broadly neutralizing HIV antibodies. *Nucleic Acids Res*481 **42** D1133 (2013).
- 482 <sup>42</sup> Raybould, M. I. J., Kovaltsuk, A., Marks, C. & Deane, C. M., CoV-AbDab: the Coronavirus Antibody
  483 Database. *bioRxiv* (2020).
- 484 <sup>43</sup> World Health Organization, WHO methods and data sources for global causes of death 2000-2016.
  485 (2016).
- 486 <sup>44</sup> Davis, C. W. *et al.*, Longitudinal Analysis of the Human B Cell Response to Ebola Virus Infection.
   487 *Cell* 177 1566 (2019).
- 488 <sup>45</sup> Arora, R. & Arnaout, R., Private Antibody Repertoires Are Public. *bioRxiv* (2020).
- 489 <sup>46</sup> Tanno, H. *et al.*, A facile technology for the high-throughput sequencing of the paired VH:VL and
- 490 TCRbeta:TCRalpha repertoires. *Sci Adv* **6** y9093 (2020).

- 491 <sup>47</sup> Goldstein, L. D. *et al.*, Massively parallel single-cell B-cell receptor sequencing enables rapid
   492 discovery of diverse antigen-reactive antibodies. *Comms Bio* 2 (2019).
- 493 <sup>48</sup> Swindells, M. B. *et al.*, abYsis: Integrated Antibody Sequence and Structure–Management, Analysis,
- 494 and Prediction. J Mol Biol **429** 356 (2017).
- 495 <sup>49</sup> Yusim, K. et al., HIV Molecular Immunology 2017. (Los Alamos National Laboratory, Theoretical
- 496 Biology and Biophysics, Los Alamos, New Mexico., 2018).
- 497 <sup>50</sup> Giudicelli, V. *et al.*, IMGT/LIGM-DB, the IMGT comprehensive database of immunoglobulin and T
- 498 cell receptor nucleotide sequences. *Nucleic Acids Res* **34** D781 (2006).
- 499 <sup>51</sup> Amid, C. *et al.*, The European Nucleotide Archive in 2019. *Nucleic Acids Res* (2019).
- 500 <sup>52</sup> Leaman, R. & Lu, Z., TaggerOne: joint named entity recognition and normalization with semi-
- 501 Markov Models. *Bioinformatics* **32** 2839 (2016).
- 502 <sup>53</sup> Larkin, M. A. *et al.*, Clustal W and Clustal X version 2.0. *Bioinformatics* **23** 2947 (2007).
- 503 <sup>54</sup> Eguchi, Y., PHYLIP-GUI-Tool (PHYGUI): adapting the functions of the graphical user interface for
- 504 the PHYLIP package. *J Biomed Sci Eng* **4** 90 (2011).
- 505 <sup>55</sup> Letunic, I. & Bork, P., Interactive Tree Of Life (iTOL) v4: recent updates and new developments.
  506 *Nucleic Acids Res* 47 W256 (2019).
- 507 <sup>56</sup> Schritt, D. *et al.*, Repertoire Builder: high-throughput structural modeling of B and T cell receptors.
   508 *Mol Syst Des Eng* **4** 761 (2019).
- 509

# 510 Acknowledgements

511 This study was supported by the National Natural Science Foundation of China (NSFC) (31771479) (Z. Z.), 512 NSFC Projects of International Cooperation and Exchanges of NSFC (61661146004), and the Local 513 Innovative and Research Teams Project of Guangdong Pearl River Talents Program (2017BT01S131). We 514 thank Jun Chen from MOE Laboratory of Biosystems Homeostasis & Protection and Innovation Center for

515 Cell Signaling Network, College of Life Sciences, Zhejiang University for the valuable comments, discussions,516 and suggestions.

## 517 Author Contributions

- 518 Y. Z., H. Z., Y. Z., C. L., X. Y., Y. Z., Y. C., Y. Z., J. W., C. W., C. M., and S. C. collected the datasets and
- 519 performed the bioinformatics analyses. Y. Z. and Q. X. developed and implemented the RAPID platform. M.
- 520 W., Q. W., H. Tang., W. X., and J. G. collected the samples and conducted the biological experiments. M. W.
- 521 and S. G. prepared the libraries and ran the Illumina sequencing. C. L. coordinated the project. X. Y. and Z. Z.
- 522 conceived the project. All authors were involved in the manuscript writing.

# 523 Competing of interests

524 The authors declared no competing financial interests.



**Supp. Fig. 1** Workflow of known antibody database construction. The first two boxes record the total number of sequences downloaded from 7 databases with Genbank and FASTA formats. Each procession on sequences is marked near arrow between intermediate results.

CAGGTTCAGCTGGTGC/	AGTCTGGAGGTGAAGTGAAGAAGCCTGGGGCCTCAGTGAAGGTC
TCCTGCAAGGCTTCTGC	TTACACCTTTAGCAACTATGGTATCACCTGGGTGCGACAGGCCCC
CGGACAAGGGCTTGAC	TGGCTGGGATGGATCAGCGCTTACAATGATAACACATACTATGCA
CAGAAGCTCCAGGGCA	GACTCACCATGACCACAGACACATCCACGAGCACAGCCTACATG
GAGCTGAGGAGCCTGA	GATCTGACGACACGGCCGTTTATTACTGTGCGAGAGATTACAGTG
CGCACCCCCCGGGAGG	CTACCTCCAGCACTGGGGGGGGGGGGGCACCCTGATCACCGTCTCCT
CAG	
eave empty will search.	the placeholder directly
Variable Region	curs

### b

Select an antige	n/disease	$\checkmark$
		Search
ee text search:		

**Supp. Fig. 2** Sequence and text search functions of RAPID. (a) The function of nucleotide and amino acid sequences search for both variable region and CDR3. (b) Known antibody search based on text such as antigen/disease and source id.

а



**Supp. Fig. 3 Basic information of public clones.** The number of public clone changed with the number of samples (a) and donors (b). The Y-axis were logarithmically converted with base 2. (c) The number of reads and clones from different sources. Sources of clones were defined based on types of B cell including naïve, memory, plasma, PBMCs, other (particular antigen-specific B cells), and unknown. The Y-axis were logarithmically converted with base 10.



**Supp. Fig. 4** The number of V genes for public clones shared by different number of donors. The number of donors of public clones is discrete.



Supp. Fig. 5 The substitution frequencies of J genes with the same CDR3aa among different donors. The darker the color, the higher the substitution frequency.



**Supp. Fig. 6** Identity of variable regions from FR1 to FR3 between therapeutic antibody and public clones. The X-axis means the divergence to germline reference and the Y-axis means the sequences identity. Different V genes are filled in different colors. Titles for subfigures separated by forward slash include inn id of therapeutic antibody, the number of samples and donors with such CDR3aa, and target of therapeutic antibody. Dots of therapeutic antibodies are larger than that of clones identified from Rep-seq datasets. Sub-figures are sorted according to the number of donors.



Germline divergence (%)



RMSD of clones with the same CDR3aa as Camrelizumab between species

	Species	Camrelizumab	Human	Mus musculus	Rattus norvegicus	6
	Camrelizumab					_
	Human	0.496				
	Mus musculus	1.036	0.972			
	Rattus norvegicus	0.545	0.490	0.872		_
h						
D	0 10 20	30 4	0 50	60 7	70 80 9	0 100
Humar Mus Rat	MQIPQAPWPVVWAVLQLGWRP .WVR.V.SFTS.QS .WVQ.V.SFTS.QS	WFLDSPDRPWNPPTFSP .L.EV.NG. RSL.Y. .L.EVLNK.R.L	ALLVVTEGDNATF W.T.S.A	CSFSNTSESFVLNWYR L.W.DLM.N. W.DLK.	MSPSNQTDKLAAFPEDRSQP LE.QCNGL LE.QCNGY	GQDCRFRVTQL V.A.QII. VR.A.QIV.
	110     120       PNGRDFHMSVVRARRNDSGTYI	130 14	Contraction 150	160 17 TAHPSPSPRPAGOFOTL RY	70 180 1 VVGVVGGLLGSLVLLVWV .I.IMSA.V.IPV.L.A.A .IVIMSV.V.IPV.L.A.A	200 LAVICSRAARG F. TSMSE AF. TGMSE
	210 220 TIGARRTGOPLKEDPSAVPVF AR.GSKDDT.E.A.P ARE.G.KED.P.AHA.A.P	230 22 SVDYGELDFQWREKTPEPI A.EGL A.EG	10 250 PVPCVPEQTEYAT ТАН.	260 23	70 280 ADGPRSAQPLRPEDGHCSWP .LQGPR.P.H. .QGPR.P.H.	1 L
	Ident	ity of protein seque	ences for PD	CD1 between spo	ecies	_
	Species	Human	Mus mu	sculus	Rattus norvegicus	
	Human					
	Mus musculus	59.31%				
	Rattus norvegicus	60.69%	86.4	6%		

**Supp. Fig. 7** Structure and sequence similarity of anti-PD1 clones and PDCD1 from different species. (a) Structure similarity of anti-PD1 clones. The upper panel stands for structures of Camrelizumab and clones from Human, Mus musculus, and Rattus norvegicus. Table in the lower panel records the RMSD of structures between paired species. (b) Identity of protein sequences for PDCD1 from human, Mus musculus, and Rattus norvegicus. The upper panel shows the multiple sequences alignment for them and the lower panel shows the sequences' identity.



**Supp. Fig. 8** Maturation pathway of clones with the same CDR3aa of MT658807. Variable region sequences with the same CDR3aa as MT658807 were extracted and compared with MT658807 and its' germline reference. The germline reference was chosen as root of phylogenetic tree and MT658807 is marked by arrow. The cluster map contains four layers including similarity of sequences (the sequences extracted from the same donor were marked with the same color), V gene family, isotype, and somatic hypermutation rate from inner to outer.



Supp. Fig. 9 Overlap of public clones shared by other species.