

1 **B cell receptor repertoire analysis in six immune-mediated**
2 **diseases**

3
4
5 RJM Bashford-Rogers^{1,2*}, L Bergamaschi^{1,3}, EF McKinney^{1,3}, DC Pombal^{1,3}, F
6 Mescia^{1,3}, JC Lee^{1,3}, DC Thomas¹, SM Flint^{1,5}, P Kellam⁴, DRW Jayne¹, PA Lyons^{1,3},
7 KGC Smith^{1,3*}

8
9 ¹Department of Medicine, University of Cambridge, Cambridge, United Kingdom.

10 ²Wellcome Centre for Human Genetics, University of Oxford, Oxford, United
11 Kingdom.

12 ³Cambridge Institute for Therapeutic Immunology and Infectious Disease, University
13 of Cambridge, Cambridge, United Kingdom.

14 ⁴ Department of Medicine, Division of Infectious Diseases, Imperial College, London,
15 United Kingdom.

16 ⁵Current address: Immunoinflammation Therapy Area Unit, GlaxoSmithKline,
17 Stevenage, United Kingdom.

18
19 ***Correspondence:** rbr1@well.ox.ac.uk, kgcs2@medschl.cam.ac.uk
20
21
22
23

24 **Introductory Paragraph**

25
26 B cells are important in the pathogenesis of many, and perhaps all, immune-mediated
27 diseases (IMDs). Each B cell expresses a single B cell receptor (BCR)¹, with the
28 diverse range of BCRs expressed by an individual's total B cell population being
29 termed the "BCR repertoire". Our understanding of the BCR repertoire in the context
30 of IMDs is incomplete, and defining this could reveal new insights into pathogenesis
31 and therapy. We therefore compared the BCR repertoire in systemic lupus
32 erythematosus (SLE), ANCA-associated vasculitis (AAV), Crohn's disease (CD),
33 Behçet's disease (BD), eosinophilic granulomatosis with polyangiitis (EGPA) and IgA
34 vasculitis (IgAV), analysing BCR clonality, and immunoglobulin heavy chain gene
35 (*IGHV*) and, in particular, isotype usage. An IgA-dominated increased clonality in SLE
36 and CD, together with skewed *IGHV* gene usage in these and other diseases,
37 suggested a microbial contribution to pathogenesis. Different immunosuppressive
38 treatment had specific and distinct impacts on the repertoire; B cells persisting after
39 rituximab were predominately isotype-switched and clonally expanded, the inverse of
40 those persisting after mycophenolate mofetil. A comparative analysis of the BCR
41 repertoire in immune-mediated disease reveals a complex B cell architecture,
42 providing a platform for understanding pathological mechanisms and designing
43 treatment strategies.
44
45
46
47
48
49
50

51

52 Main Text

53

54 Immunoglobulin gene recombination during B cell development in the bone marrow
55 (or fetal liver)² forms the “naïve” repertoire, which is modified by the
56 removal/suppression of self-reactive B cells to reduce the chance of autoimmune
57 disease³ (although 20-40% of B cells remain autoreactive⁴). Further repertoire
58 diversification occurs after B cells respond to antigen. Many undergo “isotype
59 switching” where stepwise DNA deletion and recombination from IgM generates
60 “downstream” isotypes (IgG1/2/3/4, IgA1/2, IgD and IgE) which confer distinct
61 functional characteristics and roles in disease^{5,6}. Isotype delineation is thus vital for a
62 full analysis of the BCR repertoire. Further BCR diversification occurs in specialized
63 germinal centers (GCs) – where V gene somatic hypermutation (SHM) may enhance
64 BCR affinity and specificity⁷. This post-antigenic diversification of B cell clones is
65 tempered by tolerance “checkpoints” to reduce the risk of autoimmunity⁸. The
66 peripheral BCR repertoire is thus a composite of both the naïve repertoire and that
67 generated by antigenic encounter.

68

69 BCR repertoire features have been correlated with both microbial interactions and
70 IMDs, with specific IGHV regions recognizing commensal and/or pathogenic microbes
71 or being associated with IMDs (Table S1). We analysed the BCR repertoire in 209
72 individuals across six IMDs (Tables S2, Extended Data Figure 1a), comparing (i) IMDs
73 characterized by autoantibody responses against either single dominant (AAV) or
74 multiple (SLE) autoantigens, (ii) those not thought to be autoimmune (CD, BD), and
75 (iii) those with incomplete evidence of B cell involvement or autoimmunity: EGPA
76 (formerly Churg-Strauss syndrome) and IgAV (formerly Henoch-Schönlein purpura)
77 (disease descriptions, Supplementary discussion file 1).

78

79 We developed a method to barcode, amplify, and sequence BCR repertoires from
80 RNA encoding the antigen-binding (IgH (VDJ)) and constant regions of the BCR heavy
81 chain, facilitating isotype class/subclass analysis while allowing quantitation of clone
82 frequency and correction of PCR/sequencing error (Extended Data Figure 1b)⁹. We
83 then analyzed the BCR repertoire in sorted B cells from 19 healthy controls
84 (Supplementary discussion file 2, Extended Data Figure 1-2) to develop methods to
85 control for the impact of age and differential cellular RNA content (Methods, Extended
86 Data Figure 2-3a-c, Table S4). We define the “normalized” isotype usages
87 representing the percentage of unique VDJ sequences per isotype, thus counting each
88 B cell’s contribution to the repertoire only once.

89

90 Comparative studies in IMDs have often been confounded by differences in disease
91 duration, activity and treatment. We therefore specifically recruited patients with
92 objective evidence of active disease and had not yet commenced treatment (although
93 stable doses of low-level therapy known not to affect repertoire were permitted;
94 methods, Supplementary discussion file 1). The majority were newly diagnosed. In all
95 patients the number of B cells sampled was less than the number of unique BCR
96 sequences detected (Table S3). We compared isotype use in repertoires from
97 unseparated peripheral blood mononuclear cells (PBMC) in healthy controls and IMD
98 patients (Figure 1a-b, Extended Data Figure 3d). Compared to health, IgA was over-
99 represented in all diseases except AAV and EGPA, particularly so in SLE and CD.
100 This corresponded with increased serum IgA most pronounced in SLE (Figure 1c). IgE

101 was raised in SLE, CD and, in particular, EGPA (Figure 1b, Extended Data Figure
102 3d,e), which also exhibited elevated IgG3. Isotype usage in AAV was similar to healthy
103 controls. There is therefore marked variation in isotype use in IMD, with IgA the
104 unexpected dominant isotype in diseases such as SLE and BD.

105
106 BCR repertoire diversity is driven in part by differential use of *IGHV* genes, as well as
107 non-template additions/deletions. Some individual genes, and *IGHV* subgroups
108 (defined by structural similarity¹⁰), preferentially bind microbial antigens and/or have
109 been associated with autoimmunity (Table S1). We examined *IGHV* gene frequency
110 in naïve and antigen-experienced B cells across IMDs (Figure 1d, Extended Data
111 Figure 4a, Supplemental data files 2-3). *IGHV4* family genes were increased in CD,
112 SLE and EGPA, as was *IGHV6-1*. Interestingly, *IGHV4-34* binds both autoantigens¹¹
113 and commensal bacteria¹², and has been associated with SLE¹³. Our data extends the
114 SLE association of *IGHV4-34* (and its 9G4 idiotype) to EGPA and CD (Extended Data
115 Figure 4b). Both *IGHV6-1* and *IGHV4-59* have been associated with autoreactivity
116 (Table S4). *IGHV* gene associations were seen in both the predominantly “naïve” and
117 “post-antigenic” compartments, and in both non-expanded and expanded clones
118 (Extended Data Figure 4a), raising the possibility they are not purely a consequence
119 of selected expansion after disease development (except in CD where *IGHV*
120 differences were predominantly “pre-antigenic”). V1 family genes were over-
121 represented in IMDs, particularly CD and BD. The most striking association was of BD
122 with *IGHV1-46*, -3 and -69, all previously associated with infection, in both the naive
123 and post-antigenic repertoires. Reduced representation of *IGHV* genes is also seen in
124 some diseases, reflecting either a proportional reduction due to increased frequency
125 of other *IGHV* genes or real disease associations. Levels of SHM did not vary between
126 diseases (Extended Data Figure 4c).

127
128 Increased length of complementary determining region 3 (CDR3) of the BCR is
129 associated with antibody polyreactivity and autoimmunity¹⁴. Building on previous
130 work⁹, we found an association between CDR3 length and *IGHV* gene use in healthy
131 individuals (Extended Data Figure 2, Extended Data Figure 4d). In disease, increased
132 CDR3 length was found in SLE (IgG and IgA) and CD (unswitched B cells) (Extended
133 Data Figure 4c).

134
135 B cell clones are defined by sharing a unique VDJ rearrangement, and can be
136 characterised by size (clonal expansion) and diversification (due to SHM and isotype
137 switching). Using a “clone sampling” repertoire visualization method (Supplementary
138 discussion file 4, Extended Data Figure 5), we found no differences between healthy
139 controls and AAV or IgAV, reduced clonality in BD, but increased clonal expansion
140 and complexity in CD, EGPA and SLE (Supplementary data file 5,6). We extended
141 this analysis by determining the Clonal Expansion Index (a measure of “unevenness”
142 of the number of RNA molecules per unique VDJ region sequence via the Gini index¹⁵)
143 and Clonal Diversification Index (measuring the unevenness of unique VDJ region
144 sequences per clone) (see Methods, Figures 1e-g, Extended Data Figure 6). CD
145 patients had increased clonal expansion and diversification across many isotypes,
146 particularly IgA, IgG and IgM. SLE showed a similar pattern, though with increased
147 clonality primarily in unswitched cells and with greater variation between patients, as
148 did EGPA, but with IgE predominant. Differences in maximum clone size were
149 consistent with these data (Extended Data Figure 6c,d,7a). In contrast, patients with
150 active AAV or IgAV showed no gross difference in clonal expansion or diversification,

151 and in BD both were reduced compared to controls. We then used a multivariate
152 comparison to assess “clonal normality” (see Methods), and found significant
153 dissimilarity between the repertoires of CD, EGPA and SLE patients, compared to
154 healthy, AAV, and BD patients (Extended Data Figure 7b), reinforcing the concept that
155 while some diseases are associated with broad abnormalities of the BCR repertoire,
156 others are comparatively normal.

157
158 Class-switch recombination (CSR) is a deletional DNA recombination process, so the
159 order of constant regions on the chromosome defines the possible isotypes to which
160 any given B cell can switch (Figure 2a, Extended Data Figure 7c-d). Progression of
161 CSR between each possible constant region (‘switch events’) may be assessed by
162 quantifying the frequency of unique VDJ regions sharing two isotypes (suggesting their
163 common clonal origin; Figure 2b) after normalising for read depth (Extended Data
164 Figure 7e,f). The class-switch types detectable in this analysis are reduced by the
165 isotype ambiguity between IgA1/2 and IgG1/2 in the isotype-specific sequencing, and
166 by alternative splicing of IgD from IgM-containing transcripts (Extended Data Figure
167 8a-b). We confirmed reported switch event frequencies in healthy individuals¹⁶ (Figure
168 2c). Switching differences between isotypes in IMDs usually corresponded with
169 differences in isotype usage (Figure 2d). All switching was reduced in AAV and BD,
170 and that between IgM and IgG in IgAV. In SLE and CD, increased IgA representation
171 and IgA and IgE switching was seen. The elevated isotype switching in CD appeared
172 independent of isotype frequency. In EGPA increased switching to IgE from all
173 isotypes was striking (Figure 2d), particularly IgG3 - perhaps secondary to increased
174 IgG3 frequency (Extended Data Figure 9a,b). This first systematic analysis of isotype
175 switching in IMD reveals disease-specific increases that contribute to isotype profiles.
176 Some of these, such as the prominence of IgG3/IgE in EGPA, and reduced switching
177 in AAV and BD, were unexpected and may be relevant to disease pathogenesis.

178
179 This BCR repertoire analysis supports suggestions in the literature^{17,18} that, like the
180 mouse¹⁹, human IgE clones might usually arise from clonally diversified memory cells
181 of precursor isotypes. Consistent with this IgE+ peripheral blood B cells are commonly
182 plasmablasts (Extended Data Figure 2), are unusually likely to share a clonal origin
183 with non-IgE cells, and have fewer IgE closest clonal relatives (Figure 2e-f, Extended
184 Data Figure 9c,d, Supplementary discussion file 2). IgE also commonly arises from
185 multiple independent switch events in large clones (Figure 2g).

186
187 Different immunosuppressive regimens have different impacts on the B cell
188 compartment, and these can correlate with clinical efficacy²⁰. We investigated the
189 effect of treatment on the BCR repertoire in SLE and AAV, taking repeat samples at 3
190 or 12 months after diagnosis. Most patients were treated with rituximab (RTX, B cell-
191 depleting anti-CD20 monoclonal antibody), or mycophenolate mofetil (MMF, inosine
192 5'-monophosphate dehydrogenase inhibitor precursor predominantly impacting
193 proliferating cells²¹). These regimens were standardized but not formally protocolized,
194 reflecting real-world practice based on international guidelines, and were
195 accompanied by similar steroid and subsequent maintenance therapy (Supplementary
196 discussion file 1, Methods), allowing their effect on BCR repertoire to be compared.

197
198 MMF and RTX had markedly different impacts on repertoire (Figure 3a-e, Extended
199 Data Figure 10a-c). MMF therapy resulted in an increased proportion of IgM+/D+ B
200 cells and concomitantly reduced isotype-switched B cell number and clonality, with

201 relative preservation of both SHM⁺ and SHM⁻ IgM clones compared to switched clones.
202 This could be consistent with a short half-life for switched but not IgM memory B cells
203 in humans (as seen in the mouse²²), adding to the ongoing debate on this topic²³⁻²⁴.
204 Conversely, after RTX, circulating B cell numbers were low²⁵ but persisting cells were
205 largely isotype-switched and clonally expanded, predominantly IgA in AAV and IgG1/2
206 in SLE. Larger studies are required to determine if these repertoire changes associate
207 with disease subsets, pathogenic clonal persistence, and/or treatment efficacy.
208 Nonetheless, this suggests that B cell receptor repertoire impact might inform the
209 design of therapeutic strategies (e.g. the ability of MMF to reduce class-switched
210 clones might suggest efficacy in preventing relapse following RTX therapy).

211
212 Clonal persistence 3 months after therapy was observed in >90% patients, with the
213 isotype of persistent clones differing between therapies (Figures 3f, Extended Data
214 Figure 10d-e, Methods). In AAV, reduced persistence of isotype-switched clones was
215 associated with reduced ANCA titre (Figure 3g). Persistent clones could expand,
216 undergo SHM and isotype switch despite continuing therapy (Figure 3h). By
217 considering the time between the last MMF or RTX dose and sample collection, we
218 could analyse repertoire “recovery”. After MMF, the isotype-switched population
219 reached healthy levels after approximately one year (Figure 3i). In contrast, the slow
220 reconstitution of IgD⁺/M⁺ unmutated cells after RTX is consistent with the known
221 kinetics of B cell recovery after such depletion²⁶.

222
223 This study reveals profound variation in many aspects of the BCR repertoire across
224 IMDs, both at diagnosis and after therapy. Many of the disease-associated changes
225 in isotype use, in particular, have not been previously described. The B cell receptor
226 repertoire changes in these diseases illustrate deficiencies in our understanding of
227 disease pathogenesis (Supplementary discussion file 2). SLE, CD and EGPA
228 exhibited abnormal isotype-specific clonal expansion/diversity and *IGHV* gene use,
229 such broad repertoire dysregulation being consistent with their associations with
230 multiple antibodies. Increased IgA was expected in an intestinal disease like CD, but
231 not in SLE where IgG is implicated in pathogenesis and intestinal inflammation is not
232 prominent²⁷. These observations suggests unanticipated commonality in
233 pathogenesis of SLE, CD and EGPA, suggesting they might share unknown drivers,
234 perhaps within the mucosal microbiome given known *IGHV* affinities for microbial
235 antigens^{11-13,28}. EGPA also displayed IgG3 expansion and disproportionate switching
236 to IgE. The IgE association was expected²⁹, but whether expanded IgG3 is important
237 to EGPA pathogenesis remains uncertain. IgAV associated with increased IgA and
238 mucosal involvement, but showed no evidence of IgA clonal expansion or abnormal
239 *IGHV* gene usage, consistent with distinct pathogenesis from CD. It is also possible to
240 have severe active autoimmune disease, such as AAV, without detectable B cell
241 receptor repertoire changes – the pathogenic anti-MPO or anti-PR3 clones
242 presumably being too infrequent to skew PBMC-level repertoire analysis. Finally BD
243 shows a marked increase in *IGHV1-46*, *IGHV1-69* and *IGHV1-3*, all of which bind to
244 both microbial antigens and autoantigens (Table S4), enhancing speculation that
245 infection might drive disease³⁰. Future expanded repertoire studies with, for example,
246 comparison to the microbiome or determination of the antigenic specificity of
247 expanded clones, would be illuminating. Altogether, this comprehensive analysis of
248 the B cell receptor repertoire across diseases reveals a complex architecture, which
249 may provide a platform for better understanding pathological mechanisms and
250 designing treatment strategies.

251 **References**

252

- 253 1 Nossal, G. J. V. & Lederberg, J. Antibody production by single cells. *Nature* **181**, 1419-
254 1420 (1958).
- 255 2 Lydyard, P. M., Whelan, A. & Fanger, M. W. Instant Notes Series; Instant Notes in
256 Immunology. i-x, 1-318 (2000).
- 257 3 Nemazee, D. Mechanisms of central tolerance for B cells. *Nat Rev Immunol* **17**, 281-
258 294, doi:10.1038/nri.2017.19 (2017).
- 259 4 Wardemann, H. *et al.* Predominant autoantibody production by early human B cell
260 precursors. *Science* **301**, 1374-1377, doi:10.1126/science.1086907 (2003).
- 261 5 Stavnezer, J. & Schrader, C. E. IgH chain class switch recombination: mechanism and
262 regulation. *J Immunol* **193**, 5370-5378, doi:10.4049/jimmunol.1401849 (2014).
- 263 6 Stavnezer, J., Guikema, J. E. & Schrader, C. E. Mechanism and regulation of class
264 switch recombination. *Annu Rev Immunol* **26**, 261-292,
265 doi:10.1146/annurev.immunol.26.021607.090248 (2008).
- 266 7 De Silva, N. S. & Klein, U. Dynamics of B cells in germinal centres. *Nat Rev Immunol*
267 **15**, 137-148, doi:10.1038/nri3804 (2015).
- 268 8 Giltiay, N. V., Chappell, C. P. & Clark, E. A. B-cell selection and the development of
269 autoantibodies. *Arthritis Res Ther* **14 Suppl 4**, S1, doi:10.1186/ar3918 (2012).
- 270 9 Petrova, V. N. *et al.* Combined Influence of B-Cell Receptor Rearrangement and
271 Somatic Hypermutation on B-Cell Class-Switch Fate in Health and in Chronic
272 Lymphocytic Leukemia. *Frontiers in Immunology* **9**, doi:10.3389/fimmu.2018.01784
273 (2018).
- 274 10 Matsuda, F. *et al.* The complete nucleotide sequence of the human immunoglobulin
275 heavy chain variable region locus. *J Exp Med* **188**, 2151-2162 (1998).
- 276 11 Pascual, V. *et al.* Nucleotide sequence analysis of the V regions of two IgM cold
277 agglutinins. Evidence that the VH4-21 gene segment is responsible for the major
278 cross-reactive idiotype. *J Immunol* **146**, 4385-4391 (1991).
- 279 12 Schickel, J. N. *et al.* Self-reactive VH4-34-expressing IgG B cells recognize commensal
280 bacteria. *J Exp Med* **214**, 1991-2003, doi:10.1084/jem.20160201 (2017).
- 281 13 Tipton, C. M. *et al.* Diversity, cellular origin and autoreactivity of antibody-secreting
282 cell population expansions in acute systemic lupus erythematosus. *Nat Immunol* **16**,
283 755-765, doi:10.1038/ni.3175 (2015).
- 284 14 Meffre, E. *et al.* Immunoglobulin heavy chain expression shapes the B cell receptor
285 repertoire in human B cell development. *J Clin Invest* **108**, 879-886,
286 doi:10.1172/JCI13051 (2001).
- 287 15 Bashford-Rogers, R. J. M. *et al.* Network properties derived from deep sequencing of
288 human B-cell receptor repertoires delineate B-cell populations. *Genome Res* **23**,
289 1874-1884, doi:10.1101/gr.154815.113 (2013).
- 290 16 Horns, F. *et al.* Lineage tracing of human B cells reveals the in vivo landscape of
291 human antibody class switching. *Elife* **5**, doi:10.7554/eLife.16578 (2016).
- 292 17 Saunders, S. P., Ma, E. G. M., Aranda, C. J. & Curotto de Lafaille, M. A. Non-classical B
293 Cell Memory of Allergic IgE Responses. *Front Immunol* **10**, 715,
294 doi:10.3389/fimmu.2019.00715 (2019).
- 295 18 Croote, D., Darmanis, S., Nadeau, K. C. & Quake, S. R. High-affinity allergen-specific
296 human antibodies cloned from single IgE B cell transcriptomes. *Science* **362**, 1306-
297 1309, doi:10.1126/science.aau2599 (2018).

298 19 He, J. S. *et al.* IgG1 memory B cells keep the memory of IgE responses. *Nat Commun*
299 **8**, 641, doi:10.1038/s41467-017-00723-0 (2017).

300 20 Jayne, D. R., Gaskin, G., Pusey, C. D. & Lockwood, C. M. ANCA and predicting relapse
301 in systemic vasculitis. *QJM* **88**, 127-133 (1995).

302 21 Karnell, J. L. *et al.* Mycophenolic acid differentially impacts B cell function depending
303 on the stage of differentiation. *J Immunol* **187**, 3603-3612,
304 doi:10.4049/jimmunol.1003319 (2011).

305 22 Tarlinton, D. & Good-Jacobson, K. Diversity among memory B cells: origin,
306 consequences, and utility. *Science* **341**, 1205-1211, doi:10.1126/science.1241146
307 (2013).

308 23 Seifert, M. & Kuppers, R. Human memory B cells. *Leukemia* **30**, 2283-2292,
309 doi:10.1038/leu.2016.226 (2016).

310 24 Macallan, D. C. *et al.* B-cell kinetics in humans: rapid turnover of peripheral blood
311 memory cells. *Blood* **105**, 3633-3640, doi:10.1182/blood-2004-09-3740 (2005).

312 25 Mei, H. E. *et al.* Steady-state generation of mucosal IgA⁺ plasmablasts is not
313 abrogated by B-cell depletion therapy with rituximab. *Blood* **116**, 5181-5190,
314 doi:10.1182/blood-2010-01-266536 (2010).

315 26 Anolik, J. H. *et al.* Delayed memory B cell recovery in peripheral blood and lymphoid
316 tissue in systemic lupus erythematosus after B cell depletion therapy. *Arthritis*
317 *Rheum* **56**, 3044-3056, doi:10.1002/art.22810 (2007).

318 27 Villalta, D. *et al.* Anti-dsDNA antibody isotypes in systemic lupus erythematosus: IgA
319 in addition to IgG anti-dsDNA help to identify glomerulonephritis and active disease.
320 *PLoS One* **8**, e71458, doi:10.1371/journal.pone.0071458 (2013).

321 28 Bende, R. J. *et al.* Identification of a novel stereotypic IGHV4-59/IGHJ5-encoded B-
322 cell receptor subset expressed by various B-cell lymphomas with high affinity
323 rheumatoid factor activity. *Haematologica* **101**, e200-203,
324 doi:10.3324/haematol.2015.139626 (2016).

325 29 Manger, B. J. *et al.* IgE-containing circulating immune complexes in Churg-Strauss
326 vasculitis. *Scand J Immunol* **21**, 369-373 (1985).

327 30 Galeone, M., Colucci, R., D'Erme, A. M., Moretti, S. & Lotti, T. Potential Infectious
328 Etiology of Behcet's Disease. *Patholog Res Int* **2012**, 595380,
329 doi:10.1155/2012/595380 (2012).

330
331

332 Figure legends

333

334 **Figure 1. Differences in isotype, IGHV gene usages and clonality between IMDs.**

335 **a)** Heatmap of the normalized isotype usages per disease. **b)** The percentages of
336 normalized IgA1/2 and IgE BCR percentage usages per disease. **c)** IgA titre in healthy
337 individuals (n=4) and CD (n=20) and SLE (n=8) patients. **d)** Heatmap of *IGHV* gene
338 frequency and BCR subtypes in health and disease: IgM⁺D⁺SHM⁻ BCR sequences,
339 (>78% derived from naive B cells); IgM⁺D⁺SHM⁺ BCR sequences - SHM is evidence
340 of antigenic stimulation; and IgM⁻D⁻ BCR sequences, all isotype-switched and
341 therefore post-antigenic. Light and dark orange squares indicate significantly higher,
342 and light and dark blue squares lower, gene frequency in disease than health. Only
343 genes >0.1% in frequency are shown. Relative mean gene frequencies in healthy
344 individuals are indicated at the top (full heatmap in Supplemental data file 3). *IGHV*
345 genes are ordered according to amino acid similarity, indicated by the *IGHV* gene
346 amino acid similarity tree (see methods). **e)** Explanations of clonality measures and
347 network representations of BCR repertoires. **f)** Heatmap showing the (left) Clonal
348 Expansion Index and (right) Clonal Diversification Index of each isotype per disease
349 from total PBMC B cells. **g)** The (left) Clonal Expansion Index and (right) Clonal
350 Diversification Index for PBMC BCR repertoires per disease. For (a),(b),(d),(f),(g):
351 n=32, 18, 32, 12, 10, 23, 10 and 13 for healthy, AAV MP0+, AAV PR3+, EGPA, SLE,
352 CD IgAV and Behçet's patients respectively. P-values calculated by two-sided ANOVA
353 for (a)-(c) and * denotes false discovery rate (FDR) <0.05, ** <0.005, *** <0.0005,
354 where FDR determined by Šidák method. For (d),(f): Light and dark orange squares
355 indicate significantly higher, and light and dark blue squares significantly lower, isotype
356 use in disease compared to health. Boxplots show the 25th, 50th and 75th percentiles;
357 whiskers show upper and lower quartiles.

358

359 **Figure 2. Class-switching in IMDs.**

360 **a)** Schematic diagram of class-switch recombination (CSR). **b)** Relative frequencies
361 of CSR between different constant regions may be determined through the frequency
362 of unique VDJ regions expressed as two isotypes normalized for read depth. **c)**
363 Relative class-switch recombination event frequencies across healthy individuals
364 (n=32). **d)** CSR frequencies across autoimmune diseases. Each circle represents an
365 isotype class per disease, where the size is proportional to percentage of unique BCRs
366 corresponding to that isotype, coloured according to whether it is significantly higher
367 (red), lower (blue) or not different (black) to healthy individuals. Arrows indicate class-
368 switching between isotypes, where the thickness is proportional to the relative class-
369 switch recombination event frequencies for each disease, colored according to
370 whether these are significantly higher (red), lower (blue) or not different (black) to
371 healthy individuals. P-values calculated by two-sided ANOVA, FDR determined by
372 Šidák method and threshold of significance is FDR<0.05. **e)** The proportion of VDJ
373 sequences per isotype that are also observed as other isotypes (across health and
374 diseases, n=149). **f)** The proportion of VDJ sequences where closest clonal relatives
375 are also present as same isotype (across health and diseases, n=149). **g)** A
376 representative phylogenetic tree of an IgE-associated clone maintained over the
377 course of therapy from an EGPA ANCA- patient (patient 145). Colors indicate isotype
378 usage and time-point for each BCR. All nodes scaled to unitary size. For (c),(e)-(g):
379 Boxplots show the 25th, 50th and 75th percentiles; whiskers show upper and lower
380 quartiles. For (e)-(g): P-values calculated by two-sided Wilcoxon tests for (a)-(c); *
381 denotes FDR <0.05, ** <0.005, *** <0.0005, and FDR determined by Šidák method.

382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407

Figure 3. The impact of therapy on the B cell receptor repertoire.

a) Mean proportion and phenotype distribution of B cells within PBMCs in healthy controls and in AAV and SLE patients before and after therapy (B cell percentage of PBMC in brackets). Data from AAV and SLE patients is combined post-therapy. **b-d)** Percentages of **b)** unmutated IgD/IgM, mutated IgD/IgM and antigen-experienced switched BCRs (IgA1/2/IgG1/2/IgG3/IgE), **c)** Clonal expansion, **d)** Clonal Diversification Indices, and **e)** ratio of the percentage of IgM⁺SHM⁺ BCRs over class-switched BCRs of AAV and SLE patient samples at diagnosis (red, *untreated*), and patients 3-months post-treatment with MMF (blue) or RTX (green). For AAV: Untreated (n=42), MMF (n=5), RTX (n=5), and for SLE: Untreated (n=11), MMF (n=6), RTX (n=9). **f)** Isotype percentages of persistent clones between diagnosis and post-induction therapy in AAV patients: MMF (n=5) and RTX (n=6). **g)** Percentages of persistent BCRs shared between diagnosis and 3-months, or between 3- and 12-months post-induction therapy respectively, split between patients who became serum ANCA-negative after induction versus those who remain serum ANCA positive. **h)** Percentages of persistent clones expanded by >2 fold, changed <2-fold or decreased by >2-fold between diagnosis and post-induction therapy in AAV patients (n=12, 20 and 19 for 0-3-months, 0-12-months and 3-12-months respectively). **i)** Correlation between proportions of BCR types with time since last treatment of (top) MMF (n=27), and (bottom) RTX (n=26). Pearson's correlation p-values indicated. Healthy individual BCR frequencies in red. For (b)-(h): P-values calculated by two-sided ANOVA, * denotes FDR <0.05, **<0.005, ***<0.0005, FDR determined by Šidák method. Boxplots show the 25th, 50th and 75th percentiles; whiskers show extent of outliers.

408 **Materials and methods (online)**

409

410 ***Ethical approval***

411 Ethical approval for this study was obtained from the Cambridge Local Research
412 Ethics Committee (reference numbers 04/023, 08/H0306/21, 08/H0308/176) and
413 Eastern NHS Multi Research Ethics Committee (07/MRE05/44), with informed
414 consent obtained from all subjects enrolled.

415

416 ***Samples***

417

418 **Healthy participants**

419 Inclusion criteria for healthy individuals were people aged between 20-77 years, with
420 no serious co-morbidities, no direct family history of autoimmune disease, no use of
421 immunosuppressants or steroids, and no hospitalization within the last 12 months. The
422 healthy individual samples used for B cell sorting were recruited through the NIHR
423 Cambridge BioResource.

424

425 **Patients with AAV**

426 AAV patients attending or referred to the specialist vasculitis unit at Addenbrooke's
427 Hospital, Cambridge, UK, between July 2004 and June 2016 were enrolled. Active
428 disease at presentation was defined by at least 1 major or 3 minor Birmingham
429 Vasculitis Activity Score (BVAS) criteria³¹ and the clinical impression that induction
430 immunosuppression would be required. Prospective disease monitoring was
431 undertaken monthly with serial BVAS assessment³¹ and serum ANCA status
432 (Supplementary discussion file 1). 41/54 patients were sampled at diagnosis and
433 13/54 patients at disease flare as defined above. A minority of patients (11/54) had
434 received prior treatment with oral prednisolone, and 3 patients had received
435 Azathioprine within 6 months prior to sampling. Patients on low dose steroids and
436 azathioprine have been analysed separately, and their inclusion does not impact upon
437 any of the findings described in this study.

438

439 **Patients with SLE**

440 The SLE cohort comprised patients attending or referred to the Addenbrooke's
441 Hospital specialist vasculitis unit between July 2004 and June 2016 who met at least
442 four American College of Rheumatology SLE criteria³², presenting with active disease.
443 Active disease was defined as meeting all three of the following prospectively defined
444 criteria: new British Isles Lupus Assessment Group (BILAG) score A or B in any
445 system, clinical assessment of active disease by the reviewing physician and increase
446 in immunosuppressive therapy as a result. After treatment with an
447 immunosuppressant, patients were followed up monthly. Disease monitoring was
448 undertaken with serial BILAG assessment and serum ANA status. Patients' treatment
449 was at the physician's discretion, not dictated by study participation and includes
450 therapy used for induction of remission at enrolment ('induction'). 8/10 patients were
451 sampled at diagnosis and 2/10 patients at disease flare. A minority of patients (3/10)
452 had received prior treatment with oral prednisolone and/or hydroxychloroquine.

453

454 **Patients with CD**

455 Patients with active Crohn's disease were recruited from a specialist IBD clinic at
456 Addenbrooke's Hospital, before starting treatment. 22/23 patients were recruited at
457 the time of diagnosis. Diagnosis was made using standard endoscopic, histological

458 and radiological criteria³³. All patients had at least moderately active Crohn's disease
459 at enrolment as evidenced by clinical symptoms in conjunction with some or all of
460 elevated C-reactive protein, elevated fecal calprotectin, radiologically active disease
461 or endoscopically active disease. All patients were treatment naïve, with none
462 receiving immunomodulators, corticosteroids or biological therapy.

463

464 **Patients with CLL**

465 Patients with CLL were recruited from the specialist leukemia/lymphoma unit at
466 Addenbrooke's Hospital unit between January 2011 and July 2014. CLL patient
467 inclusion required the presence of at least 5×10^9 B cells/L circulating clonal B cells
468 persisting for 3 months and a characteristic phenotype (typically CD5, CD19, CD20,
469 and CD23).

470

471 **Patients with EGPA**

472 EGPA patients attending or referred to the specialist vasculitis unit at Addenbrooke's
473 Hospital, Cambridge, UK, between July 2004 and June 2016 were enrolled into the
474 present study. EGPA diagnosis was based on the history or presence of *both* asthma
475 and eosinophilia ($>1.0 \times 10^9/L$ and/or $> 10\%$ of leukocytes) *plus* at least two additional
476 features of EGPA, criteria used in the recent Phase III clinical trial "Study to Investigate
477 Mepolizumab in the Treatment of Eosinophilic Granulomatosis with Polyangiitis". 7/11
478 patients were sampled at diagnosis and 4/11 patients at disease flare. A minority of
479 patients (4/11) had received prior treatment with oral steroids (methylprednisolone or
480 prednisolone), 2/11 patients treated with azathioprine and 1/11 patients treated with
481 cyclophosphamide within 6 months of sampling.

482

483 **Patients with IgAV and Behçet's Disease**

484 IgAV patients and Behçet's disease patients were recruited from the specialist
485 vasculitis clinic at Addenbrooke's Hospital were enrolled into the present study
486 between 2005 and 2015. Clinical data recorded for Behçet's disease patients
487 comprised: (i) Basis for diagnosis i.e. orogenital mucosal ulceration, prior ocular
488 inflammation, and characteristic skin rash (erythema nodosum or pseudofolliculitis);
489 (ii) Major complications such as venous or arterial thrombosis, central nervous system
490 involvement or involvement of the pulmonary vascular system; and (iii) disease activity
491 (expert physician global assessment). 5/11 of patients had received prior treatment
492 with oral steroids (prednisolone) and 3/11 patients had been treated with azathioprine
493 within 6 months prior to sampling.

494

495 The diagnosis of IgAV was based on the American College of Rheumatology 1990
496 criteria for the classification of Henoch-Schönlein purpura³⁴ and the 2012 Revised
497 International Chapel Hill Consensus Conference Nomenclature of Vasculitides³⁵. All
498 patients had to have a biopsy-proven diagnosis of IgAV. Patient inclusion was based
499 on if they had i) severe involvement of at least 1 organ (including biopsy-proven IgAV-
500 related nephritis class 3–4; gastrointestinal involvement with haemorrhage, ischemia,
501 perforation, and/or abdominal pain unresponsive to common analgesics and lasting
502 for >24 hours; pulmonary haemorrhage, episcleritis, cardiac and central nervous
503 system involvement); and ii) other systemic autoimmune or neoplastic diseases were
504 excluded. 8/10 IgAV patients were sampled at diagnosis and 2/10 patients at disease
505 flare. 4/10 of patients had received prior treatment with oral prednisolone, 1/10 patients
506 treated with azathioprine and 1/10 patients treated with cyclophosphamide within 6
507 months of sampling.

508

509 **Cell separation, RNA extraction and antibody titres**

510 For PBMCs and CD19+ B cells: PBMCs were isolated from 110 ml of whole blood by
511 centrifugation over Ficoll. CD19+ B cells were isolated by positive selection using
512 magnetic beads as previously described³⁶. Total RNA was extracted from each
513 sample using an RNeasy mini kit (Qiagen) with quality assessed using an Agilent
514 BioAnalyser 2100 and RNA quantification performed using a NanoDrop ND-1000
515 spectrophotometer.

516

517 For flow-sorted B cell samples from Espéli et al.³⁷: Flow sorting was performed using
518 CD19-BV785, CD38-BV711, CD3-NC650, CD14-605NC, CD24-PerCP-Cy5.5, IgD-
519 FITC, CD27-PE-Cy7 and Aqua (Invitrogen) (flow protocol outlined in Extended Data
520 Figure 1), into sorting buffer (10mM Tris pH 8.0 and RiboLock RNase Inhibitor (1U/μL))
521 and frozen immediately.

522

523 Total IgA and IgE levels in patient serum were measured using a ProcartaPlex
524 immunoassay kit (ThermoFisher) using 25ul of serum from each individual and run on
525 a Luminex xMAP analyser. Raw data (MFI) were normalised to a concurrently
526 measured 7 point standard curve according to the manufacturer's instructions to return
527 an absolute quantification (pg/ml). All measured values were encompassed by the
528 standard distribution.

529

530 **Reverse transcription and amplification with barcoded primers**

531 Reverse transcription (RT) was performed in a 23uL reaction: 14ul of RT mix 1
532 (containing RNA template, 10uM reverse primer mix, 1 μL dNTP (10mM), and
533 nuclease-free water) was incubated for 5 min at 70°C. This mixture was immediately
534 transferred to ice for 1 min, and the RT mix 2 (4 μL 5x FS buffer, 1 μL DTT (0.1M), 1
535 μL SuperScript®III (Thermo Fisher)) was added and incubated at 50°C for 60 min
536 followed by 15 min inactivation at 70°C. cDNA was cleaned with Agencourt AMPure
537 XP beads and PCR amplified with V-gene multiplex primer mix (10μM each forward
538 primer) and 3' universal reverse primer (10μM) using KAPA protocol and the thermal
539 cycling conditions: 1 cycle (95°C - 5 min); 5 cycles (98°C - 5 sec; 72°C - 2 min); 5
540 cycles (65°C - 10 sec, 72°C - 2 min); 25 cycles (98°C - 20sec, 60°C - 1 min, 72°C - 2
541 min); 1 step (72°C - 10 min). Primers are provided in STAR Methods.

542

543 **Sequencing and barcode filtering**

544 Sequencing libraries were prepared using Illumina protocols and sequenced using
545 300bp paired-ended sequencing on a MiSeq (Illumina). Raw reads were filtered for
546 base quality (median Phred score >32) using QUASR
547 (<http://sourceforge.net/projects/quasr/>)³⁸. Forward and reverse reads were merged if
548 they contained identical overlapping region of >50bp, or otherwise discarded.
549 Universal barcoded regions were identified in reads and orientated to read from V-
550 primer to constant region primer. The barcoded region within each primer was
551 identified and checked for conserved bases. Primers and constant regions were
552 trimmed from each sequence, and sequences were retained only if there was >80%
553 per base sequence similarity between all sequences obtained with the same barcode,
554 otherwise discarded. The constant region allele with highest sequence similarity was
555 identified by 10-mer matching to the reference constant region genes from the IMGT
556 database³⁹, and sequences were trimmed to give only the region of the sequence
557 corresponding to the variable (VDJ) regions. Isotype usage information for each BCR

558 was retained throughout the analysis hereafter. Sequences without complete reading
559 frames and non-immunoglobulin sequences were removed and only reads with
560 significant similarity to reference IGHV and J genes from the IMGT database using
561 BLAST⁴⁰ were retained. Ig gene usages and sequence annotation were performed in
562 IMGT V-QUEST, where repertoire differences were performed by custom scripts in
563 Python.

564

565 **Accounting for age in BCR analysis**

566 Age-related BCR repertoire differences have been previously described, and this
567 could be important as immune-mediated diseases often have different ages of onset.
568 We confirmed this in both healthy controls and disease by incorporating age as a
569 covariate in repertoire analyses, as in previous studies⁴¹⁻⁴³.

570

571 As expected correction for age usually made little difference (Extended Data Figure
572 4a). Where statistical discordance between uncorrected and corrected data did occur,
573 the latter became not significant, indicating this correction is appropriately
574 conservative (i.e. correction does not create spurious statistically significant positive
575 associations). In these and most other cases, predominantly in diseases of later onset
576 (AAV, EGPA), age correction made p values less significant, indicating some
577 observed repertoire differences are driven in part by age, and underlining the
578 importance of correcting for it (Extended Data Figure 3a-d). In some cases, already
579 significant results became more so after correction – as expected many of these were
580 in SLE or CD, diseases with a younger age of onset (Extended Data Figure 3c).

581

582 **Isotype frequencies, somatic hypermutation, CDR3 lengths and IGHV gene** 583 **usages**

584 To account for the greater numbers of BCR RNA molecules per plasmablast
585 compared to other B cell subsets, we calculated two measures of isotype usage: (1)
586 the percentage of reads per isotype which does not control for differential RNA per
587 cell, thus reflecting the impact of plasmablasts/plasma cells on repertoire, and (2) the
588 normalized isotype usages, defined as the percentage unique VDJ sequences per
589 isotype, thus controlling for differential RNA per cell and reducing potential biases from
590 differential RNA per cell. We did not control for ethnicity as the majority of patients
591 (95%) in all disease groups were of northern European ancestry, with the exception of
592 SLE in which 4 patients were Asian and 5 were Caucasian. We observed only two
593 *IGHV* genes with differential frequencies between ethnicities with FDR <0.05 (Table
594 S6), and neither of these were differentially expressed between SLE and health.

595

596 Similarly, mean somatic hypermutation levels and CDR3 lengths were calculated per
597 unique VDJ region sequence to reduce potential biases from differential RNA per cell.
598 *IGHV* gene usages were determined using IMGT⁴⁴, and proportions were calculated
599 per unique VDJ region sequence. The representation of *IGHV* genes in the BCR
600 repertoire reflects their presence in the germline, the naïve repertoire and their
601 expansion after antigenic exposure. We therefore compared the frequency of *IGHV*
602 gene use in PBMC-derived BCRs identified by sequence as being enriched for naïve
603 (IgM⁺D⁺SHM⁻: >78% naïve B cells by flow cytometry) and antigen-experienced B cells
604 (including both unswitched (IgM⁺D⁺SHM⁺) and class-switched memory (IgA⁺/G⁺/E⁺)
605 subsets).

606

607

608 ***BCR repertoire generation and network analysis***

609 The network generation algorithm and network properties were calculated as in
610 Bashford-Rogers *et al.*¹⁵: each vertex represents a unique sequence, where relative
611 vertex size is proportional to the number of identical reads. Edges join vertices that
612 differ by single nucleotide non-indel differences and clusters are collections of related,
613 connected vertices. A clone (cluster) refers to a group of clonally related B cells, each
614 containing BCRs with identical CDR3 regions and *IGHV* gene usage, or differing by
615 single point mutations, such as through SHM. Each cluster is assumed to arise from
616 the same pre-B cell.

617

618 Repertoire parameters that were dependent on sequencing depth were generated by
619 subsampling each sequencing sample to a specified depth:

620

621 1) The Clonal Expansion index is a measure of “unevenness” of the number of RNA
622 molecules per unique VDJ region sequence by vertex Gini Index as defined in
623 Bashford-Rogers *et al.*¹⁵. This is calculated from the distribution of the number of
624 unique RNA molecules per vertex within subsampled BCR repertoires at specified
625 depth defined below. The mean of 100 repeats of resulting Clonal Expansion
626 indices was determined.

627

628 2) The Clonal Diversification index is a measure of the unevenness of unique VDJ
629 region sequences per clone by cluster Renyi Index as defined in Bashford-Rogers
630 *et al.*¹⁵. This is calculated from the distribution of the number of unique VDJ region
631 sequences per clone within subsampled BCR repertoires at specified depth
632 defined below. The mean of 100 repeats of resulting Clonal Diversification indices
633 was determined. Clone size distributions were also calculated from the same
634 subsamples and a mean of 100 repeats was determined.

635

636 The number of sampled unique RNA molecules (for the Clonal Expansion index) and
637 clones (for the Clonal Diversification index) per sample was: all isotypes: 3500, Ig/M
638 mutated: 600, Ig/M unmutated: 500, Class-switched: 1000, IgA1/2: 1000, IgD: 75, IgE:
639 50, IgG1/2: 500, IgG3: 100, IgM: 750. These thresholds were chosen as a balance
640 between including as many samples as possible per analysis whilst remaining as
641 representative of the total BCR repertoire in each sample.

642

643 ***BCR network sampling to preserve the overall clonal structure of visual***
644 ***representation***

645 We developed network sampling methods to obtain a graphical representation of a
646 network that preserves the overall clonal architecture. The rationale for this
647 development, and for the selection of the Clone Sampling method, is discussed in
648 detail in Supplementary discussion file 4. Briefly, a fixed number of clones were
649 subsampled, and a network generated from all BCRs from these clones from a given
650 sample. Subsampling was performed 1000 times, and the sample that contained a
651 maximum clone size closest to the median of all subsamples was chosen to generate
652 a visual representation of the BCR repertoire.

653

654 ***Global measure of BCR repertoire***

655 To define a global measure of the “normality” or otherwise of the BCR repertoire, we
656 combined three main BCR features (isotype frequency, clonal expansion index, clonal

657 diversification index) using a multivariate MANOVA comparison between disease
658 groups using age as a covariate.

659

660 ***Class-switching event analyses***

661

662 **Relative class-switch event frequency** was the frequency of unique VDJ regions
663 expressed as two isotypes (i.e. from more than one B cell, where one has undergone
664 class-switch recombination). This was determined as proportion of unique BCRs
665 present as both isotypes IgX and IgY within a random subsample of 8000 BCRs,
666 where the mean of 1000 repeats was generated (Extended Data Figure 7e). This
667 provides information on the frequency of BCRs observed associated with any two
668 isotypes (class-switching events) while accounting for total read depth, but not
669 accounting for differences in the relative frequencies of BCRs per isotype.

670

671 The **per-isotype normalized class-switch event frequencies** determines frequency
672 of unique VDJ regions expressed as two isotypes whilst normalizing for differences in
673 isotype frequencies. To account for differences in isotype proportions, BCRs from
674 each isotype were randomly subsampled to a fixed depth of 100 BCRs, and the
675 proportion of unique VDJ sequences present between each pair of isotypes was
676 counted (Extended Data Figure 9a). The mean of 1000 repeats was generated.

677

678 **Clonal overlap between time points during therapy**

679 The identification of persistent clones was performed using MRDARCY⁴⁵. Clonal
680 overlap frequencies between samples, including the quantification of persistent
681 clones, was determined through subsampling each repertoire to a fixed depth of 2000
682 unique BCRs and determining the proportion of overlapping clones. The mean of 1000
683 repeats was generated.

684

685 Although quantitative conclusions are difficult as a blood draw samples such a small
686 proportion of peripheral B cells, the clonal overlap estimate between timepoints is, as
687 expected, significantly lower than that from technical BCR sequencing repeats from
688 the same RNA samples and higher than the overlap between unrelated patient
689 samples (Extended Data Figure 10d).

690

691 ***Phylogenetic analysis***

692 Phylogenetic trees from AAV and EGPA patients were derived from all clusters
693 containing at least one BCR sequence across multiple time points using the
694 MRDARCY pipeline⁴⁶. Alignments were performed using Mafft⁴⁷ and maximum
695 parsimony trees fitted using Paup^{*48}. The *IGHV* gene amino acid similarity tree was
696 generated using an alignment of reference *IGHV* genes from IMGT using Mafft⁴⁷ and
697 a maximum parsimony tree was fitted using Paup^{*48}.

698

699 ***Statistical methods***

700 Statistical differences between disease groups were performed using ANOVA or
701 MANOVA with patient age as a covariate and correcting for multiple testing by
702 Bonferroni correction. Where patients were age-matched, Wilcoxon tests were
703 performed.

704

705 **Data access**

706 Sequencing data available from the EGA (accession numbers in Table S3).

707 **Supplementary information** is available in the online version of the paper.

708

709 **Acknowledgements** This work was supported by the Wellcome Trust (grant
710 WT106068AIA and 083650/Z/07/Z), the EU H2020 project SYSCID (grant no.
711 733100), the UK Medical Research Council (program grant MR/L019027) and the UK
712 National Institute of Health Research (NIHR) Cambridge Biomedical Research Centre.
713 We gratefully acknowledge the patients who participated in this study, and Valerie
714 Morrison, Angela Reynolds, all NIHR Cambridge BioResource staff and volunteers,
715 and the Cambridge NIHR BRC Cell Phenotyping Hub (particularly Anna Petrunkina
716 Harrison, Natalia Savinykh Yarkoni, Esther Perez, Simon McCallum, and Chris
717 Bowman). We thank Federico Alberici, Nuru Noor and other members of the
718 Addenbrooke's Vasculitis and Gastroenterology services, Norberto Escudero Urquijo
719 for discussions about network subsampling, Plamena Naydenova and Giulia
720 Manferrari. We are grateful to John Todd and David Tarlinton for reviewing the
721 manuscript.

722

723 **Author contributions** R.B.R. and K.G.C.S. planned the study. R.B.R performed
724 BCR amplification, and analysed sequencing data. F.M. analysed clinical data and
725 L.B., D.C.P. and S.M.F. performed immunophenotyping. E.F.N., J.C.L., D.C.T.,
726 S.M.F., D.R.W.J, and P.A.L. contributed to sample collection and clinical data
727 generation, and P.K. contributed to sample processing. R.B.R., P.A.L., E.F.N., J.C.L.,
728 D.C.T. and K.G.C.S. provided intellectual contributions to analyses. R.B.R. and
729 K.G.C.S. wrote the manuscript. All authors edited manuscript.

730

731 **Author information** Reprints and permissions information is available at
732 www.nature.com/reprints. Correspondance and requests for materials should be
733 addressed to RBR (rbr1@well.ox.ac.uk) or KS (kgcs2@medschl.cam.ac.uk).
734 Competing financial interests: Rachael Bashford-Rogers, Paul Kellam and Ken Smith
735 are all named on a patent associated with the methodologies in this paper. Shaun Flint
736 is a current employee of GlaxoSmithKline, and holds shares in GlaxoSmithKline.
737 Rachael Bashford-Rogers is a consultant for Imperial College London and
738 VHSquared. Paul Kellam is an employee and holder of shares in Kymab Ltd. David
739 Jayne is a recipient of a research grant from Roche/Genetech.

740

741 **References**

742

- 743 31 Stone, J. H. *et al.* A disease-specific activity index for Wegener's granulomatosis:
744 modification of the Birmingham Vasculitis Activity Score. International Network for
745 the Study of the Systemic Vasculitides (INSSYS). *Arthritis Rheum* **44**, 912-920,
746 doi:10.1002/1529-0131(200104)44:4<912::AID-ANR148>3.0.CO;2-5 (2001).
- 747 32 Tan, E. M. *et al.* The 1982 revised criteria for the classification of systemic lupus
748 erythematosus. *Arthritis Rheum* **25**, 1271-1277 (1982).
- 749 33 Silverberg, M. S. *et al.* Toward an integrated clinical, molecular and serological
750 classification of inflammatory bowel disease: report of a Working Party of the 2005
751 Montreal World Congress of Gastroenterology. *Can J Gastroenterol* **19 Suppl A**, 5A-
752 36A (2005).
- 753 34 Mills, J. A. *et al.* The American College of Rheumatology 1990 criteria for the
754 classification of Henoch-Schonlein purpura. *Arthritis Rheum* **33**, 1114-1121 (1990).
- 755 35 Jennette, J. C. *et al.* 2012 revised International Chapel Hill Consensus Conference
756 Nomenclature of Vasculitides. *Arthritis Rheum* **65**, 1-11, doi:10.1002/art.37715
757 (2013).
- 758 36 Lyons, P. A. *et al.* Microarray analysis of human leucocyte subsets: the advantages of
759 positive selection and rapid purification. *BMC Genomics* **8**, 64, doi:10.1186/1471-
760 2164-8-64 (2007).
- 761 37 Espeli, M. *et al.* FcγRIIb differentially regulates pre-immune and germinal
762 center B cell tolerance in mouse and human. *Nat Commun* **10**, 1970,
763 doi:10.1038/s41467-019-09434-0 (2019).
- 764 38 Watson, S. J. *et al.* Viral population analysis and minority-variant detection using
765 short read next-generation sequencing. *Philos Trans R Soc Lond B Biol Sci* **368**,
766 20120205, doi:10.1098/rstb.2012.0205 (2013).
- 767 39 Lefranc, M. P. IMGT unique numbering for the variable (V), constant (C), and groove
768 (G) domains of IG, TR, MH, IgSF, and MhSF. *Cold Spring Harb Protoc* **2011**, 633-642,
769 doi:10.1101/pdb.ip85 (2011).
- 770 40 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
771 search tool. *J Mol Biol* **215**, 403-410, doi:10.1016/S0022-2836(05)80360-2 (1990).
- 772 41 Davydov, A. N. *et al.* Comparative Analysis of B-Cell Receptor Repertoires Induced by
773 Live Yellow Fever Vaccine in Young and Middle-Age Donors. *Front Immunol* **9**, 2309,
774 doi:10.3389/fimmu.2018.02309 (2018).
- 775 42 Marioni, R. E. *et al.* Genetic Stratification to Identify Risk Groups for Alzheimer's
776 Disease. *J Alzheimers Dis* **57**, 275-283, doi:10.3233/JAD-161070 (2017).
- 777 43 Ellis, J. A., Panagiotopoulos, S., Akdeniz, A., Jerums, G. & Harrap, S. B. Androgenic
778 correlates of genetic variation in the gene encoding 5α-reductase type 1. *J Hum
779 Genet* **50**, 534-537, doi:10.1007/s10038-005-0289-x (2005).
- 780 44 Giudicelli, V., Chaume, D. & Lefranc, M. P. IMGT/V-QUEST, an integrated software
781 program for immunoglobulin and T cell receptor V-J and V-D-J rearrangement
782 analysis. *Nucleic Acids Res* **32**, W435-440, doi:10.1093/nar/gkh412 (2004).
- 783 45 Bashford-Rogers, R. J. *et al.* Eye on the B-ALL: B-cell receptor repertoires reveal
784 persistence of numerous B-lymphoblastic leukemia subclones from diagnosis to
785 relapse. *Leukemia* **30**, 2312-2321, doi:10.1038/leu.2016.142 (2016).

786 46 Bashford-Rogers, R. J. M. *et al.* Eye on the B-ALL: B-cell receptor repertoires reveal
787 persistence of numerous B-lymphoblastic leukemia subclones from diagnosis to
788 relapse. *Leukemia*, doi:10.1038/leu.2016.142 (2016).

789 47 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
790 improvements in performance and usability. *Molecular biology and evolution* **30**,
791 772-780, doi:10.1093/molbev/mst010 (2013).

792 48 Wilgenbusch, J. C. & Swofford, D. Inferring evolutionary trees with PAUP*. *Current*
793 *protocols in bioinformatics / editorial board, Andreas D. Baxevanis ... [et al.] Chapter*
794 **6**, Unit 6 4, doi:10.1002/0471250953.bi0604s00 (2003).

795

796 **Extended Data Figure Legends:**

797

798 **Extended Data Figure 1. Overview of BCR repertoire strategy.**

799 **a)** Schematic diagram of the B cell receptor repertoire analysis strategy. **b)** Schematic
800 diagram of the B cell receptor sequencing strategy. In the reverse transcription (RT)
801 step, the primer anneals to the constant region of the B cell receptor mRNA to
802 generate cDNA with a random 12 nucleotide barcode. This barcode can be
803 computationally used to reduce PCR amplification biases after sequencing. The
804 product is cleaned and PCR amplified using multiple primers binding to the FR1 region
805 of the IgH genes along with a universal sequence complementary to the end of the RT
806 primer. **c)** Gating strategy to flow sort B cell subsets from healthy donor PBMCs.

807

808 **Extended Data Figure 2. Impact of B cell subset and age on repertoire.**

809 **a)** The isotype usage frequencies from BCR sequencing data from sorted naïve B cells
810 (CD19+IgD+CD27-), CD19+CD27-IgD- B cells, IgD+ memory/B1/MZ B cells
811 (CD19+CD27+IgD+) and IgD- memory B cells (CD19+CD27-IgD-CD38-) and
812 plasmablasts (CD19+CD27+IgD-CD38+) from 19 healthy individuals. **b)** (left) The
813 mean CDR3 lengths and (right) mean SHM per BCR from healthy individual cell-sorted
814 B cell populations (n=19). **c)** The plasmablast frequency per patient in peripheral blood
815 at enrolment as a percentage of CD19+ B cells. **d)** Distribution of patient ages within
816 this study, grouped by disease. **e-g)** PBMC BCR repertoire correlations with age in
817 healthy individuals for **g)** the mean number of somatic hypermutation per BCR per bp,
818 **f)** the percentages of BCRs per isotype; **g)** percentage sizes of the largest cluster per
819 sample. For b-c: P-values calculated by two-sided by ANOVA and * denotes FDR
820 <0.05, ** <0.005, *** <0.0005, where FDR determined by Šidák method. For e-g: P-
821 values calculated by two-sided by Wilcoxon test. * denotes p-values <0.05, ** <0.005,
822 *** <0.0005, NS otherwise, with raw p-values provided in Table S4. Boxplots show the
823 25th, 50th and 75th percentiles; whiskers show upper and lower quartiles.

824

825 **Extended Data Figure 3. Repertoire changes with age and isotope usages with**
826 **disease.**

827 **a)** Correlation of p-values obtained using age as a covariate versus excluding age in
828 the analysis across 178 BCR features (calculated by two-sided by ANOVA), and **b)**
829 analyses where statistical significance was discordant (i.e. below significance
830 threshold without accounting for age and above significance threshold while using age
831 as a covariate, or vice versa, purple points from (a)). **c)** Analyses where statistically
832 significant p-values were decreased by >1.5 fold after using age as a covariate. Grey
833 dotted lines in (a) indicate the threshold of significance after account for multiple
834 testing (FDR<0.05, determined by Šidák method). **d)** The percentages of normalized
835 isotype usages (unique VDJ sequences per isotype) for PBMC BCR repertoires per
836 disease. **e)** Normalized IgE immunoglobulin constant region transcript levels between
837 disease groups from transcriptomic data. n=58, 23, 33, 13, 10, 8, 11 and 37 for healthy,
838 AAV MP0+, AAV PR3+, Behçet's, CD, EGPA, IgAV and SLE patients respectively.
839 **f)** IgE titre between healthy individuals (n=4) and EGPA patients (n=5). P-values
840 calculated by two-sided by ANOVA for (d)-(e) and * denotes FDR <0.05, ** <0.005, ***
841 <0.0005, where FDR determined by Šidák method. Boxplots show the 25th, 50th and
842 75th percentiles; whiskers show upper and lower quartiles.

843

844 **Extended Data Figure 4. Changes in IGHV gene usage with disease.**

845 Changes in *IGHV* gene usage between unexpanded and expanded clones. **a)**
846 Heatmap of each *IGHV* gene frequency difference between healthy individuals and
847 each autoimmune disease patient group within BCRs from IgM⁺D⁺ or isotype-switched
848 (IgA/IgG/E) BCR from unexpanded clones (containing <3 unique BCRs) or expanded
849 clones (3 or more unique BCRs per clone). Only genes >0.1% in frequency are shown.
850 *IGHV* genes are ordered according to amino acid similarity as in Figure 2. **b)** *IGHV4-*
851 *34* BCR frequencies with autoreactive AVY & NHS motifs compared between healthy
852 individuals and disease groups, separated by BCR type: IgM⁺D⁺SHM⁻ BCR
853 sequences, IgM⁺D⁺SHM⁺ BCR sequences and IgM⁻D⁻ BCR sequences (defined in (a)).
854 **c)** Heatmaps showing the (top) mean SHM per BCR and (bottom) relative mean CDR3
855 lengths mean SHM per BCR per isotype per disease from total PBMC B cells. **d)** The
856 distribution of the mean CDR3 lengths per *IGHV* gene in healthy individuals (n=32).
857 Each point represents a mean CDR3 length for an individual for (left) unmutated IgD/M
858 BCRs and (right) class-switched BCRs. Instances where *IGHV* genes represented by
859 fewer than 10 BCRs in an individual are excluded. For (a)-(d): n=32, 18, 32, 12, 10,
860 23, 10 and 13 for healthy, AAV MP0+, AAV PR3+, EGPA, SLE, CD IgAV and Behçet's
861 patients respectively. P-values calculated by two-sided ANOVA. Orange squares
862 indicate significantly higher, and blue squares significantly lower, corresponding gene
863 frequency between healthy individuals and disease. FDR determined by Šidák
864 method.

865

866 **Extended Data Figure 5. Network subsampling methods for preserving** 867 **repertoire structure.**

868 **a)** Schematic diagram of the cluster-vertex migration in the CC algorithm. **b)** Maximum
869 cluster sizes between true (unsampled) networks and sub-sampled networks of 2000
870 clones by the tree subsampling methods. **c)** Comparison of representative networks
871 from each patient group at diagnosis. The sample patient samples are represented
872 across the three sampling methods. Each vertex represents a unique sequence,
873 where relative vertex size is proportional to the number of identical reads. Edges join
874 vertices that differ by single nucleotide non-indel differences and clusters are
875 collections of related, connected vertices. Networks are comprised of a subsample of
876 2000 clones using the corresponding subsampling method. Each vertex is
877 represented by a pie chart indicating the percentage of each isotype, where blue =
878 IgD/M, red = IgA1/2, yellow = IgG1/2, green = IgG3, and grey = IgE.

879

880 **Extended Data Figure 6. BCR repertoire clonality between diseases.**

881 **a)** Boxplots of the Clonal Expansion Index and **b)** the Clonal Diversification Index for
882 PBMC BCR repertoires per disease. **c)** Plots of the percentage of clones per sample
883 per disease greater than clone size, C. Clone size is defined as the number of unique
884 VDJ sequences that are clonally related. For each disease, the mean percentage is
885 indicated by the dark blue line, and the upper and lower interquartile ranges indicated
886 by the light blue areas. Overlaid in grey is the equivalent for healthy individuals.
887 Differences in read depth were accounted for by subsampling 5000 clones from each
888 repertoire and determining a mean of 20 repeats. As a disease comparison, we show
889 the distribution for CLL. **d)** Boxplots of the percentage of clones larger than 10, 20, 30,
890 40, or 50 unique VDJs per disease. Differences in reads depth were accounted for by
891 performing subsamples 5000 clones and determining a mean of 20 repeats. For
892 (a),(b),(d): P-values calculated by two-sided by ANOVA for (a),(b),(d). * denotes FDR
893 <0.05, ** <0.005, *** <0.0005, and FDR determined by Šidák method. n=32, 18, 32,
894 12, 10, 23, 10 and 13 for healthy, AAV MP0+, AAV PR3+, EGPA, SLE, CD IgAV and

895 Behçet's patients respectively. Boxplots show the 25th, 50th and 75th percentiles;
896 whiskers show upper and lower quartiles.

897

898 **Extended Data Figure 7. BCR repertoire similarity between diseases and class-**
899 **switch recombination estimation.**

900 **a)** The maximum clone sizes (as a percentage of unique VDJ sequences of a given
901 isotype in largest clone divided by the total of unique BCRs of that isotype) for PBMC
902 BCR repertoires per disease across isotypes. **b)** Global repertoire dissimilarity
903 measures between disease groups. Heatmap showing the global repertoire
904 dissimilarity measures between disease groups based on the combination of three
905 main BCR features (isotype frequency, clonal expansion index, clonal diversification
906 index) and determining joint differences between groups (MANOVA test using disease
907 group and age as covariables). The light and dark orange squares indicate significant
908 differences between corresponding disease groups (false discovery rate (FDR) < 0.05
909 and 0.005 respectively, where FDR determined by Šidák method. n=32, 18, 32, 12,
910 10, 23, 10 and 13 for healthy, AAV MP0+, AAV PR3+, EGPA, SLE, CD IgAV and
911 Behçet's patients respectively. **c)** The sequence of B cell isotype expression is defined
912 by the order of constant regions on the chromosome, where the possible class-
913 switching is depicted by the arrows between constant regions. **d)** Schematic diagram
914 of class-switch types detectable from the sequencing data due to the ambiguity of
915 isotype between IgA1/2 and IgG1/2 in the isotype-specific sequencing and splicing of
916 IgD from IgM-containing transcripts. Possible class-switching events are represented
917 by the arrows between constant regions. **e)** Multiple unique RNA sequences with
918 identical antigen binding regions (V-D-J) but different constant regions represent
919 instances of class switching. **f)** Schematic diagram of sub-sampling of BCR repertoires
920 to generate the relative class-switch event frequency. This is the frequency of unique
921 VDJ regions expressed as two isotypes (i.e. from more than one B cell, where one has
922 undergone CSR), and determined as proportion of unique BCRs present as both
923 isotypes IgX and IgY within a random subsample of 8000 BCRs, where the mean of
924 1000 repeats was generated. This provides information on the frequency of BCRs
925 observed associated with any two isotypes (class-switching events) while accounting
926 for total read depth, but not accounting for differences in the relative frequencies of
927 BCRs per isotype. For (a): n=32, 18, 32, 12, 10, 23, 10 and 13 for healthy, AAV MP0+,
928 AAV PR3+, EGPA, SLE, CD IgAV and Behçet's patients respectively. P-values
929 calculated by two-sided ANOVA, * denotes FDR <0.05, ** <0.005, *** <0.0005, and
930 FDR determined by Šidák method. Boxplots show the 25th, 50th and 75th percentiles;
931 whiskers show upper and lower quartiles.

932

933 **Extended Data Figure 8. Class-switch recombination estimation differences**
934 **between diseases.**

935 **a)** Boxplots of the proportion of class-switch events between isotypes for each
936 autoimmune disease. Boxplots show the 25th, 50th and 75th percentiles; whiskers show
937 upper and lower quartiles. **b)** Boxplots of the proportion of class-switch events
938 between autoimmune diseases across isotypes for PBMC BCR repertoires via
939 subsampling total repertoire. P-values calculated by two-sided ANOVA and * denotes
940 FDR <0.05, ** <0.005, *** <0.0005, where FDR was determined by the Šidák method.
941 n=32, 18, 32, 12, 10, 23, 10 and 13 for healthy, AAV MP0+, AAV PR3+, EGPA, SLE,
942 CD IgAV and Behçet's patients respectively. Boxplots show the 25th, 50th and 75th
943 percentiles; whiskers show upper and lower quartiles. **c)** Phylogenetic trees of
944 representative clonal expansions from patients demonstrating class-switch

945 recombination events. Each vertex is a unique BCR sequence and is represented by
946 a pie chart indicating the percentage of each isotype, where blue = IgD/M, red =
947 IgA1/2, yellow = IgG1/2, green = IgG3, and grey = IgE. Branch lengths are estimated
948 by maximum parsimony, and the BCRs with the lowest number of somatic
949 hypermutations are indicated (denoted “BCRs closest to germline”).

950

951 **Extended Data Figure 9. Normalised class-switch recombination estimation**
952 **differences between diseases and IgE clonal features.**

953 **a)** Schematic diagram of sub-sampling of BCR repertoires to generate the per-isotype
954 normalized class-switch event frequencies, defined as the frequency of unique VDJ
955 regions expressed as two isotypes whilst normalizing for differences in isotype
956 frequencies. To account for differences in isotype proportions, BCRs from each
957 isotype were randomly subsampled to a fixed depth of 200 reads, and the proportion
958 of unique VDJ sequences present between each pair of isotypes was counted. The
959 mean of 1000 repeats was generated. **b)** Boxplots of the proportion of the per-isotype
960 normalized class-switch event frequencies between isotypes for each autoimmune
961 disease. P-values calculated by two-sided by ANOVA and * denotes FDR <0.05, **
962 <0.005, *** <0.0005, where FDR determined by Šidák method. n=32, 18, 32, 12, 10,
963 23, 10 and 13 for healthy, AAV MP0+, AAV PR3+, EGPA, SLE, CD IgAV and Behçet’s
964 patients respectively. Boxplots show the 25th, 50th and 75th percentiles; whiskers
965 show upper and lower quartiles. **c)** Boxplots of the mean cluster sizes per patient per
966 isotype as a percentage of BCRs per isotype, comparing IgE-associated clones with
967 non-IgE-associated clones for each disease. **d)** The proportion of VDJ sequences per
968 isotype that are observed also as other isotypes for each disease. P-values calculated
969 by two-sided Wilcoxon tests and * denotes FDR <0.05, ** <0.005, *** <0.0005, where
970 FDR determined by Šidák method. n=32, 18, 32, 12, 10, 23, 10 and 13 for healthy,
971 AAV MP0+, AAV PR3+, EGPA, SLE, CD IgAV and Behçet’s patients respectively.
972 Boxplots show the 25th, 50th and 75th percentiles; whiskers show upper and lower
973 quartiles.

974

975 **Extended Data Figure 10. Impact of therapy on BCR repertoire.**

976 The **a)** percentages of BCRs per isotype, **b)** mean SHM pre BCR per isotype, and **c)**
977 clonal expansion indices of AAV and SLE patient samples taken at diagnosis (red,
978 *untreated*), and patients post 3-months induction therapies with MMF (blue) or RTX
979 (green). For AAV, the patients per group were: Untreated (n=42), MMF (n=5), RTX
980 (n=5), and for SLE, the patients per group were: Untreated (n=11), MMF (n=6), RTX
981 (n=9). **d)** The percentage of BCRs shared between diagnosis and 3 or 12 months post-
982 induction therapy AAV samples (blue), BCRs shared between repertoires from the
983 same RNA tube (red), and BCRs shared between samples from unrelated patient
984 samples. Zero overlap was found between unrelated samples, whereas significantly
985 higher overlap between BCRs shared between repertoires from the same RNA tube
986 compared to BCRs shared between diagnosis and 3 or 12 months post-induction
987 therapy AAV samples. This suggests that the overlap measurements yield realistic
988 and normalized values at this sampling depth. **e)** The percentages of persistent BCRs
989 shared between diagnosis and 3 months post induction therapy, split between patients
990 receiving different therapies. P-values calculated by two-sided Wilcoxon tests.
991 Boxplots show the 25th, 50th and 75th percentiles; whiskers show upper and lower
992 quartiles.

993

994