

1 **Multimodal repertoire analysis unveils B cell biology in health and immune-mediated diseases**

2

3 Mineto Ota^{1,2*}, Masahiro Nakano^{1,3}, Yasuo Nagafuchi^{1,2}, Satomi Kobayashi¹, Hiroaki Hatano^{1,4},
4 Ryochi Yoshida¹, Yuko Akutsu¹, Takahiro Itamiya¹, Atsushi Matsuo⁵, Yumi Tsuchida¹, Hirofumi
5 Shoda¹, Kazuhiko Yamamoto^{1,3}, Kazuyoshi Ishigaki⁴, Tomohisa Okamura², Keishi Fujio^{1*}

6

7 ¹Department of Allergy and Rheumatology, Graduate School of Medicine, The University of Tokyo,
8 Tokyo, Japan

9 ²Department of Functional Genomics and Immunological Diseases, Graduate School of Medicine, The
10 University of Tokyo, Tokyo, Japan

11 ³Laboratory for Autoimmune Diseases, RIKEN Center for Integrative Medical Sciences, Yokohama,
12 Kanagawa, Japan

13 ⁴Laboratory for Human Immunogenetics, RIKEN Center for Integrative Medical Sciences, Yokohama,
14 Kanagawa, Japan

15 ⁵Research Division, Chugai Pharmaceutical Co., Ltd., Kamakura, Kanagawa, Japan

16 * Corresponding authors. Email: otam@m.u-tokyo.ac.jp (M.O.), fujiok-int@h.u-tokyo.ac.jp (K.F.)

17

18 **Abstract**

19 Despite involvement of B cells in the pathogenesis of immune-mediated diseases, biological
20 mechanisms underlying their function are scarcely understood. To overcome this gap, comprehensive
21 analysis of the B cell repertoire is essential. Here, we cataloged and investigated the repertoire of five
22 B cell subsets from 595 cases under immune-mediated diseases and health. CDR-H3 length among
23 naïve B cells was shortened among autoimmune diseases in an interferon signature-dependent manner.
24 VDJ gene usage was skewed especially in plasmablasts and unswitched-memory B cells of systemic
25 lupus erythematosus patients with frequent usage of VDJ genes used mainly in naïve B cells and not
26 unswitched-memory B cells of healthy controls. We developed a scoring system for this skewing, and
27 it correlated with peripheral helper T cell transcriptomic signatures and disease activity and decreased
28 after belimumab treatment. Moreover, genetic association analysis identified three molecules possibly
29 involved in somatic hyper-mutation processes in humans. Our multimodal repertoire analysis brings
30 new insights to B cell biology.
31

32 **Main**

33 B cells play critical roles in adaptive immunity, and their importance in immune-mediated disease
34 (IMD) pathogenesis is suggested by the overlap of genetic risk variants of IMDs to regulatory regions
35 in B cells^{1, 2} and by the clinical efficacy of B cell-targeted therapies in IMDs^{3, 4}. Thus, deeper
36 understanding of B cell biology could clarify IMD pathology.

37 To increase their diversity, B cell receptors (BCRs) are created by intrinsic genome-editing. In the
38 bone marrow, one V, D and J gene for one clonotype is selected (VDJ recombination), and random
39 deletions and insertions occur in the junctional regions. In the germinal center (GC), class switch
40 recombination and loading of somatic hypermutation (SHM) further increases BCR diversity. These
41 processes generate a vast variety of BCRs in an individual, a collection known as the “BCR repertoire”.
42 Heavy-chain complementarity-determining region 3 (CDR-H3) is the most variable region and central
43 to antigen specificity of BCRs.

44 For securing prompt defensive responses against diverse foreign pathogens, humans have evolved
45 a distinct immune system compared to other organisms⁵. Genetic adaptation in human populations has
46 resulted in strong immune responses against pathogens, but it was a trade-off against susceptibility to
47 autoimmunity⁵. The B cell immune system in humans is complex and has been elusive. It is now clear
48 that genetic determinants of individual difference among BCR repertoires as well as analysis of actual
49 differences in BCR repertoires among IMDs are essential for deeper understanding of B cell biology.

50 IMDs consist of autoimmune and autoinflammatory diseases⁶. Autoimmune diseases are
51 characterized by two prominent features: the emergence of autoantibodies and hyperactivity of
52 interferons (IFNs)⁶, although the relationship between these features remains elusive. For creation of
53 memory B cells and antibody secreting cells, affinity maturation through GC is a critical step. On the
54 other hand, previous studies revealed that those cell types can be differentiated in a GC-independent
55 manner⁷, a process termed “extrafollicular” maturation. Pathogenic autoantibodies are generated

56 through B cell extrafollicular reactions in systemic lupus erythematosus (SLE) and flares in SLE is
57 characterized by the expansion of naive-derived activated effector B cells of extrafollicular phenotype⁸,
58 ⁹. The disease relevance of these features of abnormal B cell maturation in IMDs could be better
59 understood if placed in the context of a large-scale dataset combined with other features.

60 We recently created a gene expression and genome sequence database containing a large variety of
61 immune cells from 10 IMD patients as well as healthy volunteers². In this follow-up study, we focused
62 on BCR features of 5 sorted B cell subsets that could be identified from RNA-seq data. We also
63 examined their associations with transcriptomic, genetic, and clinical information. This largest-ever
64 functional genomics dataset of immune cells from IMD patients enabled us to comprehensively
65 explore BCR repertoire abnormalities in IMDs. We developed a scoring system of repertoire
66 abnormalities in IMDs based on VDJ gene usage, the “Repertoire Naïveness (RN) Score”. The RN
67 score showed significant association with the transcriptomes of specific T cell subsets and disease
68 activity. Also, genetic association analysis of >500 plasmablast samples suggested previously
69 unknown biology in nucleotide substitution mechanisms in SHM. Our multi-omics approach to a
70 patient-derived, large-scale dataset generated new insights into B cell biology in both health and
71 disease states.

72

73 **Results**

74 **Profiling of BCR repertoires of IMD patients and healthy volunteers**

75 We mapped BCR regions from RNA-seq reads of 5 sorted CD19⁺ B cell subsets: CD27⁻IgD⁺ naïve B
76 cells, CD27⁺IgD⁺ unswitched memory B cells (USM B), CD27⁺IgD⁻ switched memory B cells (SM
77 B), CD27⁻IgD⁻ double negative B cells (DN B) and IgD⁻CD27⁺⁺CD38⁺ plasmablasts (Fig. 1a). We
78 extended our cohort from the previous report² and samples as follows: 136 patients with SLE, 90
79 patients with systemic sclerosis (SSc), 85 patients with idiopathic inflammatory myopathy (IIM), 147

80 patients with other IMDs and 137 healthy controls (HC), giving a total of 595 cases (Fig. 1a). For 22
81 SLE patients, samples were also obtained after a half-year treatment with belimumab, a biologic
82 targeting B-cell activating factor (BAFF). After mapping and quality control (Methods), 5.8 million
83 productive unique CDR-H3 clonotypes were identified from 2,893 samples. On average, 2,801 distinct
84 productive clonotypes were identified in each sample (Extended data Fig. 1a). More than 5,000 BCR-
85 derived reads were identified per million RNA-seq reads in each subset (median), and in plasmablasts
86 it reached 1.1×10^5 reads (Fig. 1b), reflecting their specialized role in secreting antibodies.

87 As naïve B cells are expected to have few SHM, the mutation frequency can reflect the error rate in
88 PCR or sequencing steps. The SHM frequency in V genes in naïve B cells was well controlled at 0.15%
89 (median), in clear contrast to 6.7% (median) in plasmablasts (Fig. 1c). The number of shared CDR-H3
90 sequences was substantially larger in comparisons of the same individuals than in the comparisons of
91 different individuals (Fig. 1d). The greatest degree of sharing was observed among the samples from
92 the same individual and the same subset at different time points, indicating the valid reproducibility of
93 our BCR repertoire dataset (Fig. 1d). Equivalent to the previous single cell BCR analysis that classified
94 more than 99% of IgG⁺ and 98% of IgA⁺ cells as memory cells¹⁰, IgG or IgA were rarely detected
95 (0.4%) in naïve B cells in our dataset (Fig. 1e). Together, our dataset captured biologically reasonable
96 characteristics of BCR repertoires and enabled deep phenotyping of repertoires among IMDs.

97 In order to evaluate the influence of the number of reads for BCR repertoire characterization, we
98 down-sampled the number of reads from plasmablast RNA-seq data as well as the BCR-seq data of
99 peripheral blood mononuclear cells from the previous study¹¹ and mapped BCR regions in the same
100 way. These analyses revealed the strong influence of the number of reads on the number of identified
101 CDR-H3 clonotypes, in clear contrast to stable quantification of CDR-H3 length, V gene usage and
102 isotype frequency even with shallow depth (Extended data Fig. 1b). Thus, in our following RNA-seq
103 based BCR repertoire analysis, we focused on these relative features, rather than analyzing whole

104 clonotypes.

105

106 **In autoimmune diseases, CDR-H3 length is shortened in naïve B cells in an interferon signal**
107 **strength-dependent manner**

108 CDR-H3 loop length is one of the classic indices for the assessment of repertoire bias¹². When we
109 compared CDR-H3 length with HC, marked shortening of CDR-H3 lengths in naïve B cells (by more
110 than 1 amino acid) and DN B cells among autoimmune diseases were observed (Fig. 2a; Extended data
111 Fig. 2a). Conversely, CDR-H3 lengths in USM B cells and plasmablasts of SLE patients were longer
112 than those in HC (Fig. 2a; Extended data Fig. 2a). Among the sections of the CDR-H3 region,
113 lengthening in plasmablasts was observed in germline-coded regions but not in junction regions, in
114 contrast to shortening in overall sections in naïve B cells (Fig. 2b; Extended data Fig. 2b).

115 Observations so far were only from productive sequences. Non-productive sequences, which refer
116 to the sequences with a frameshift or stop codon, can be used as a means of studying the pre-selection
117 repertoire¹³. In our dataset, shortening of CDR-H3 lengths among autoimmune diseases was observed
118 in non-productive sequences of naïve and memory B cells (Extended data Fig. 2c). Thus, together with
119 the observed shortening of overall CDR-H3 sections in productive sequences, short CDR-H3
120 clonotypes may be generated before selection of mature naïve B cells (e.g., during the process of VDJ
121 recombination and nucleotide insertion/deletion at the junctions). This observation is analogous to
122 previous reports of shortening of CDR3 length in pre-selection T cell receptors among type 1 diabetes
123 patients¹³. Moreover, CDR-H3 lengthening in germline regions of plasmablasts and USM B in SLE
124 may be the consequence of skewed VDJ gene usage as described in the next section.

125 CDR-H3 shortening was especially prominent in naïve B cells from diseases with high IFN
126 signatures (Fig. 2a; Extended data Fig. 2d), although CDR3 length was not as divergent among light
127 chains (Extended data Fig. 2e). To investigate further, we tested the association of gene expression and

128 CDR-H3 length in naïve B cells in each disease and subsequently performed meta-analysis in order to
129 avoid the confounding of disease effects (Methods). Genes associated with a short CDR-H3 length
130 were significantly enriched in IFN signature genes (ISG, Fig. 2c, d), and the association was observed
131 in each autoimmune disease (Extended data Fig. 2f). CDR-H3 length in naïve B cells was also
132 significantly associated with ISG expression in 4 other B cell subsets, indicating the association of
133 generalized IFN activity itself with CDR-H3 length (Extended data Fig. 2g). Strong negative
134 correlation between CDR-H3 length and IFN activity was characteristic of naïve B cells among 5 B
135 cell subsets (Fig. 2e). These results suggested a role of IFN signaling at the early stage of B cell
136 development, which is consistent with previous findings of IFN-induced abnormal B cell maturation
137 in bone marrow of SLE¹⁴.

138

139 **Skewed gene usage in autoimmune diseases suggested an abnormal repertoire maturation** 140 **process**

141 It was previously suggested that there was differential usage of VDJ genes according to B cell subsets¹⁵,
142 ¹⁶. As these reports were limited to one healthy female¹⁵ and inbred strain of mice, C57BL/6¹⁶, the
143 generalizability of this finding in a large population has been elusive. In our cohort, gene usage
144 preference among 5 B cell subsets was also observed (Extended data Fig. 3a). More surprisingly, by
145 principal components analysis (PCA) of VDJ gene usage, samples were clearly separated according to
146 cell types (Fig. 3a). The separation pattern reflected the lineage relationships among the subsets, and
147 they were arranged similarly to the flow-cytometric plots according to CD27 positivity and IgD
148 positivity (Fig. 3a). This observation is analogous to previously reported clear separation of 3 B cell
149 subsets by PCA of VDJ use in C57BL/6 mice¹⁶. Those results strongly suggested that B cell selection
150 was based on VDJ gene segments independent of the sequence of their antigen-binding regions in
151 humans.

152 It has been shown that there is skewed usage of the IGHV in IMDs¹⁷ with dependence on the stage
153 of disease¹⁸, although the sample size has been limited. Notably, in the PCA space, the gene usage
154 signature was drastically skewed in the negative direction of the CD27 axis in SLE among USM B and
155 plasmablasts (Fig. 3a; Extended data Fig. 3b). With regard to IIM and SSc, the repertoire skewed in
156 the same direction in the PCA space in USM B, but not obviously in plasmablasts (Extended data Fig.
157 3b). In case-control comparisons, many differentially used genes were identified especially in
158 autoimmune diseases (Fig. 3b). In parallel with the findings in the PCA space, VDJ gene usage
159 differences between SLE and HC in USM B, as well as plasmablasts, showed clear agreement with
160 gene usage differences between naïve B and USM B in HC (Fig. 3c). These results indicated that the
161 majority of observed differential gene usage in SLE was associated with the global abnormality in B
162 cell maturation process, rather than the consequence of specific antigens.

163 Based on the observed skewing in VDJ gene usage in autoimmune diseases toward a CD27 negative,
164 “naïve like” repertoire, we next developed a score to quantify such skewing based on its VDJ gene
165 usage. We termed this factor the “Repertoire Naïveness (RN) Score”. Briefly, the RN score is
166 calculated as the sum of the product of normalized VDJ gene usage and the size of its effect, defined
167 by naïve B versus USM B comparisons in HC (Methods, Supplementary table 1). As we defined the
168 effect sizes of each gene only using HC samples, this score can be applied to inter-disease comparison
169 without bias. Consistent with the appearance in PCA, our RN score was the highest among naïve B
170 cells followed by DN B cells, and the highest among SLE patients in USM B or plasmablasts (Extended
171 data Fig. 3c, d). To validate our approach, we adopted the same scoring system to the BCR-seq data
172 derived from sorted B cells from human tonsils in a previous study¹⁹. The RN score was clearly high
173 in naïve B cells and low in memory B cells (Extended data Fig. 3e).

174 As reported previously²⁰, the use of each VDJ gene was associated with different CDR-H3 lengths
175 (Extended data Fig 3f). The genes associated with an extended CDR-H3 were preferentially used by

176 plasmablasts in SLE (Extended data Fig. 3g). Thus, the extended CDR-H3 length in SLE in this subset
177 is associated with abnormal VDJ gene usage, at least in part.

178

179 **The repertoire abnormality in autoimmune diseases may be associated with the extrafollicular**
180 **pathway and specific subpopulations of Th1 cell**

181 Extrafollicular maturation of plasmablasts has been reported as a hallmark of SLE⁸. Next, we assessed
182 the association of the observed “naïve like” repertoire and the extrafollicular maturation. Although
183 extrafollicular maturation-derived antibody secreting cells undergo SHMs, they are expected to have
184 less SHMs^{7, 9, 21} compared to GC-derived cells. In line with this notion, the frequency of SHMs in
185 plasmablasts was significantly lower in SLE patients than in HC (Extended data Fig. 4a). Furthermore,
186 it showed a significant negative correlation with the RN score (Fig. 4a). To examine the relationship
187 in clonotype levels, we divided BCR clonotypes into bins according to the frequencies of SHMs and
188 calculated RN scores in each bin. In this comparison, we observed that clonotypes with a lower number
189 of mutations showed higher RN scores (Fig. 4b). Taken together, naïve-like gene usage in plasmablasts
190 is associated with a low load of SHM. Reflecting the accumulation of mutations with aging¹¹, SHM
191 frequencies in plasmablasts showed a significant positive correlation with aging ($p < 2 \times 10^{-16}$),
192 although the RN score was not affected by aging ($p > 0.05$, Extended data Fig. 4b). The RN score
193 distinguished SLE patients from HC more accurately than did the SHM frequency (Extended data Fig.
194 4c). Thus, although well correlated, the RN score might reflect disease-associated pathways better than
195 the SHM frequency.

196 We next compared the RN scores of 136 SLE patients focusing on the gene expression of 27 immune
197 cell subsets that were obtained from the same patients. The number of genes that showed a significant
198 correlation with the RN score was the largest in Th1 cells for both plasmablasts and USM B (Fig. 4c).
199 The top associated genes in Th1 cells included *PDCDI*, *TOX* and *ICOS*, which are hallmarks of

200 recently identified peripheral helper T cells (Tph)²² (Fig. 4d). Indeed, the Th1 genes associated with
201 RN score showed significant concordance with the direction of change in the differential expression
202 analysis between PD-1^{hi} versus PD-1^l T cells from the original report²² (Fig. 4e). Thus, the Th1 genes
203 associated with the RN score may reflect the expansion of Tph-like subpopulation among Th1 cells.
204 The Tph gene expression score in Th1 cells (Tph1 score, Methods) was high in SLE (Extended data
205 Fig. 4d), consistent with the reported expansion of Tph population²³, particularly with Th1 phenotype²⁴,
206 in SLE. The Tph1 score showed a strong correlation with the RN score in SLE, both in USM B and
207 plasmablasts (Extended data Fig. 4e). Correlation of these two scores was also prominent in USM B
208 in IIM, although the correlation in plasmablasts was lower than that in SLE. These results may reflect
209 the shared and distinct pathogenicities among autoimmune diseases. The association of Tph1 score
210 and RN score was also validated in our independent SLE cohort²⁵ (Extended data Fig. 4f). Although
211 IFN signal strength was also correlated with RN scores (Extended data Fig. 4g), this association was
212 mostly mediated by the Tph1 score (Fig. 4f). The direct effect of IFN signal strength on RN score was
213 not significant ($p = 0.21$) in the mediation analysis (Fig. 4f). This result indicated the role of IFN on
214 Tph cell expansion. The proposed function of Tph cells is that they provide help to B cells outside the
215 GC^{22, 23}. Collectively, associations of RN score with low SHM frequencies and high Tph1 score
216 support the notion that the naïve-like repertoire in IMDs is a consequence of excess extrafollicular
217 maturation.

218 The RN score of USM B correlated significantly and positively with the expression of *TBX21* in
219 USM B, which codes for T-bet (Extended data Fig .4h). As T-bet⁺ B cells reportedly have distinct
220 functions²⁶ and are associated with autoimmunity²⁷, this naïveness-associated USM B population
221 might play a role in autoimmune pathogenesis.

222

223 **Repertoire abnormalities are associated with disease activity in SLE**

224 We next tested the association of BCR repertoire abnormalities and disease activity in SLE. RN scores
225 in plasmablasts and USM B showed significant positive correlations with SLEDAI-2K, a measure of
226 SLE disease activity (Fig. 5a), and these associations were still significant after adjusting for Tph1
227 scores ($p = 1.3 \times 10^{-3}$, 1.7×10^{-3}) and ISG scores ($p = 5.3 \times 10^{-4}$, 2.3×10^{-5}), indicating the clinical
228 relevance of RN scores. In addition, RN score in plasmablasts significantly correlated with total IgG
229 concentrations in sera (Extended data Fig. 5a).

230 We next asked whether specific V, D or J gene usage was associated with disease activity. In SLE,
231 among the 5 B cell subsets, IGHV4-34 usage in USM B showed by far the most significant positive
232 association with SLEDAI-2K ($p = 3.9 \times 10^{-10}$, Fig. 5b, Extended data Fig. 5b). Although the expansion
233 of antibody secreting cells with IGHV4-34 in acute SLE has been reported⁹, the clinical relevance of
234 its usage in USM B has not been clarified. IGHV4-34 is a V gene that is predominantly used in naïve
235 B cells (Extended data Fig. 3a) and its usage explained 28% of the variance of RN score in USM B
236 (Extended data Fig. 5c). Mediation analysis indicated that over one-half of the association between
237 RN scores in USM B and SLEDAI-2K was mediated by IGHV4-34 usage. This contrasted with the
238 independent direct association of RN scores in plasmablasts with SLEDAI-2K (Extended data Fig. 5d).
239 Thus, repertoire naïveness in SLE may be associated with disease activity partly via high usage of
240 IGHV4-34 in USM B.

241 It has been reported that IGHV4-34 has AVY and NHS motifs that are associated with
242 autoimmunity²⁸. Mutations in these motifs abrogate its autoreactivity and only such clones without
243 autoreactivity are selected for maturation under healthy conditions^{28,29}. IgG B cells without mutations
244 in these motifs from IRAK4- and MYD88-deficient patients demonstrated reactivity to autoantigens
245 as well as binding to commensal bacterial antigens²⁹. In our dataset, NHS and AVY motifs were
246 significantly more conserved without mutations in SLE USM B compared to HC as well as in

247 plasmablasts (Fig. 5c). Thus, IGHV4-34 clonotypes in SLE might feature binding to autoantigens.

248

249 **Belimumab attenuated abnormal repertoire in SLE**

250 For SLE, the BAFF inhibitor belimumab is one of the two biologics that has been approved by the
251 United States Food and Drug Administration³. In SLE, BAFF is induced by type I IFN stimulation and
252 supports the survival and differentiation of autoreactive B cells^{30, 31}. To assess the effect of BCR
253 modifications, we tested the BCR repertoire before and after treatment with belimumab in 22 cases of
254 SLE. Belimumab effectively reduced SLEDAI-2K and increased complement levels in our cohort (Fig.
255 5d). In order to assess the overall trend, we projected cases before and after belimumab treatment on
256 the same PCA space as Fig. 3a. Belimumab treatment significantly changed the repertoire of USM B
257 against the naiveness-associated direction, but not in the other subsets (Fig. 5e, f). Belimumab
258 significantly reduced the RN scores in USM B, although we detected no effects on Tph1 scores or IFN
259 scores (Fig. 5d). As to IGHV4-34, belimumab significantly reduced the fraction of nonmutated
260 AVY/NHS motifs in USM B and plasmablasts (Extended data Fig. 5e). Thus, together with its major
261 effect on the USM B repertoire, belimumab also reduced the autoimmune-associated fraction of BCRs
262 in plasmablasts.

263 We next divided all cases in our cohort, including healthy controls, into 6 clusters based solely on
264 the usage of VDJ genes of USM B and plasmablasts that were utilized for RN score calculation
265 (Extended data Fig. 5f). The largest cluster (cluster 1) mostly consisted of autoimmune disease
266 patients; these subjects showed higher RN scores, Tph1 scores and IFN scores than the others. Among
267 SLE patients, subjects in this group responded better to belimumab than did others when assessed by
268 delta SLEDAI-2K after a half-year treatment ($p = 0.007$, Extended data Fig. 5g). These results
269 indicated that the BCR repertoire abnormality reflected clinical heterogeneity within and between
270 diseases.

271

272 **Altered isotype usage in SLE is partly influenced by abnormal repertoire maturation**

273 With regard to isotype frequency, significant skewing was also observed in IMDs compared with HC
274 (Extended data Fig. 5h). As the frequencies of many isotypes were altered in SLE plasmablasts, we
275 assessed their association with RN scores in plasmablasts. The association between isotype frequencies
276 and RN scores differed according to class (Extended data Fig. 5i). A substantial proportion of the
277 decrease in IgA2 frequency and the increase in IgG1 frequency in SLE could be explained by the high
278 RN score in SLE. That finding might indicate an association with the extrafollicular pathway, although
279 the increase in IgA1 frequency in SLE seemed to have no association with this pathway (Extended
280 data Fig. 5i). Thus, the increase in IgA1 frequency in SLE might be due to other pathways, such as
281 involvement of mucosal immunity as suggested previously^{17, 29}.

282

283 **Analysis of genetic association reveals molecules associated with somatic hyper-mutation biology**

284 Finally, we tested the association of BCR repertoire features with genetic polymorphisms. For VDJ
285 gene usage, we observed significant cis-regulatory effects (Fig. 6a). This result is concordant with a
286 classical twin study that showed significant heritability in V gene usage³². Significant genetic
287 association was observed only in V gene usage, and not in D and J gene usage (Extended data Fig. 6a).
288 The genetic association with V gene usage was highly concordant among 5 B cell subsets (Extended
289 data Fig. 6b), indicating that polymorphisms have effects on gene preference at the step of VDJ
290 recombination.

291 SHM is likely a random process. However, the unevenness of nucleotide substitution patterns
292 indicates complex biological processes in the substitution process³³. To assess the association of SHM
293 substitution patterns and genetic polymorphisms, we developed a new framework, i.e., SHM-QTL
294 analysis. In this analysis, we quantified nucleotide substitution pattern frequencies in plasmablasts of

295 each subject and assessed its genetic association (Methods). In this way, we identified two loci that
296 showed significant genome-wide association with SHM patterns (Fig. 6b). Of note, in the *NUGGC*
297 locus, 3 conditionally independent lead variants were identified, and they were associated differently
298 with nucleotide substitution patterns (Fig. 6c, Extended data Fig. 6c). Two of the 3 lead variants were
299 missense variants and the other one was also in tight LD ($r^2 = 0.96$) with a missense variant. These
300 missense variants were located in different exons of *NUGGC*, and 2 of the 3 variants were estimated
301 to be located on the surface of the protein (Extended data Fig. 6d). Thus, these variants could change
302 the interaction of *NUGGC* with other molecules. In the *RRM2B* locus, the lead variant was in high LD
303 with a missense variant of *RRM2B* (Fig. 6d, rs1037699, $r^2 = 0.87$), and conditioning by the missense
304 variant cancelled the association of the lead variant (Fig. 6e), thus indicating the association of this
305 missense variant. In the *RRM2B* locus, conditioning analysis revealed the other associated variant,
306 rs892503. Although this variant was not in LD with any missense variant, the residual association
307 showed nearly perfect co-localization with eQTL effects in plasmablasts for *AP002852.2* (Fig. 6e,
308 $\text{coloc}^{34} PP. H4 = 0.985$), long non-coding RNA (lncRNA) expressed almost exclusively in memory B
309 cells and plasmablasts (Extended data Fig. 6e). The identified 5 independent variants contributed to
310 substitution patterns differently (Extended data Fig. 6f), and they explained 21.1% of SHM substitution
311 diversity in total (Methods). To our knowledge, this is the first assessment of genetic effects on SHM
312 substitution patterns and highlighted three molecules that may be associated with SHM biology in
313 humans.

314

315 **Discussion**

316 In this report, we characterized BCR repertoire abnormalities in IMDs using a large cohort of 595
317 cases. Our multi-layer dataset enabled us to link the BCR repertoire to other layer information,
318 including transcriptome of multiple immune cell types, genetic polymorphisms and clinical

319 information.

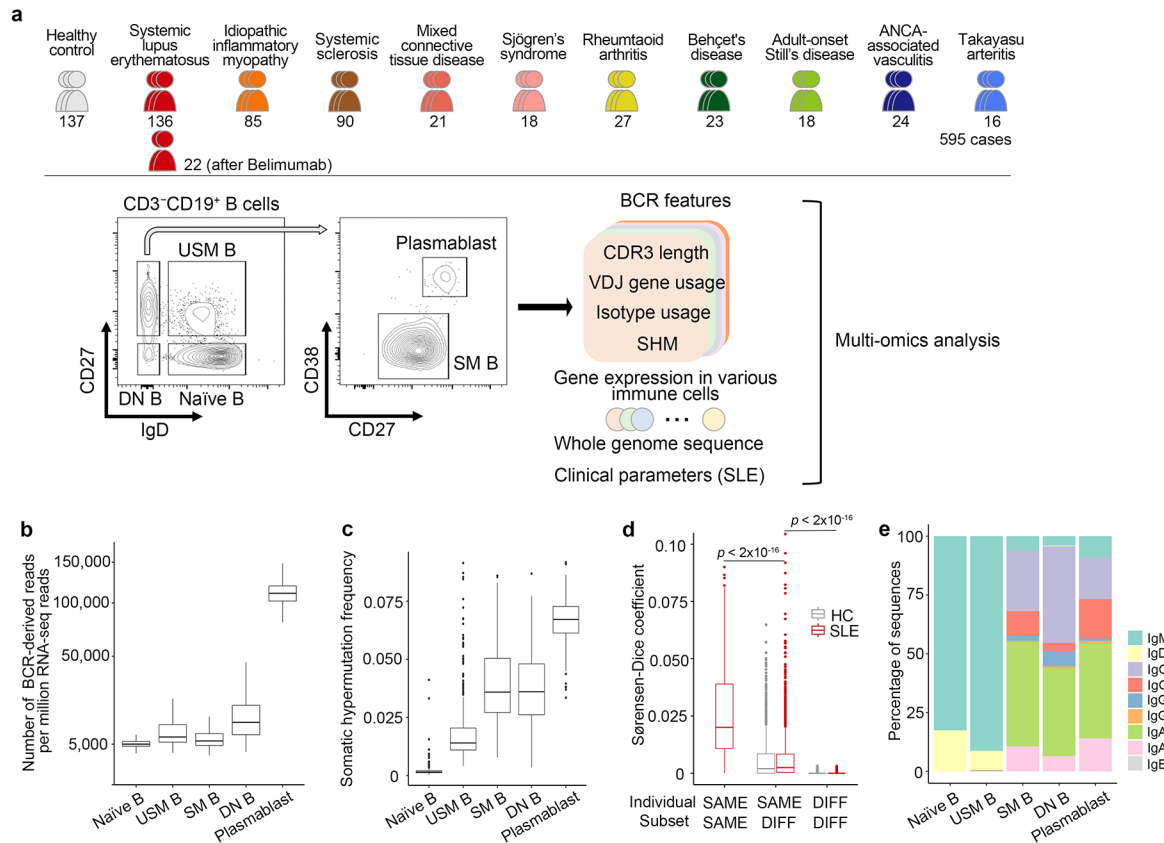
320 We identified marked shortening of CDR-H3 length in naïve B cells of autoimmune disease patients
321 in an IFN signal strength-dependent manner, as well as lengthening in plasmablasts of SLE patients.
322 Shortening of BCR CDR-H3 length in peripheral blood mononuclear cells of SLE patients had been
323 reported³⁵. Antibodies with long CDR-H3 length were associated with autoreactivity³⁶ and long CDR-
324 H3 length in class switched clonotypes in SLE was also reported¹⁷. Our results are consistent with all
325 of these seemingly contradicting reports and reinforce the importance of using sorted B cells for
326 repertoire analysis¹⁵.

327 In gene usage analysis, we observed a global pattern of skewness in specific subsets of IMDs as
328 quantified with the RN score. This score can measure extrafollicular maturation based on the
329 observation of its association with naïve-like gene usage, low SHM load and its correlation with Tph
330 signature gene expression in Th1 cells. Another possible explanation includes the defects in GC
331 selection by formation of spontaneous GCs as suggested in lupus model mice and patients^{37, 38}, which
332 might result in a similar BCR phenotype. Although the elucidation of the precise intermediate process
333 is a future problem, the observed significant association of these scores with clinical parameters and
334 their responses to clinically effective treatment suggested the disease relevance of these BCR
335 repertoire-based measures and their possible utility for clinical practice such as patient stratification.
336 Interestingly, in SLE patients, the repertoires of USM B and plasmablasts seemed to be skewed
337 similarly in response to Tph expansion. This was in contrast to the observation of similar associations
338 in USM B but less so in plasmablasts in IIM. This result may suggest the existence of a “gate-keeper”
339 for plasmablast maturation that is defective in SLE. Together with CDR-H3 length analysis, our results
340 support the notion that the break of tolerance occurs both in central and peripheral checkpoints in
341 SLE^{39, 40}. Also, our analysis supported direct and indirect associations of IFN signaling with both of
342 these checkpoints.

343 In SHM-QTL analysis, we identified 3 molecules possibly associated with SHM biology. SHM is
344 triggered by activation-induced cytidine deaminase (AID)⁴¹. Although multiple molecules are involved
345 in subsequent DNA damage response⁴², the roles of 3 identified molecules in this study has not been
346 well known. *NUGGC*, Nuclear GTPase, Germinal Center Associated, was first identified as a GTPase
347 expressed in germinal center B cells. It inhibits the function of AID in cell lines⁴³. Interestingly,
348 *NUGGC*-deficient mice reportedly have a substantial increase in mutations at G:C base pairs⁴⁴ and
349 experience altered patterns of SHM. These data, together with our findings, strongly suggest an
350 association of this molecule in DNA repair in SHM processes in humans. *RRM2B* has been reported
351 to repair DNA damage in a TP53-dependent manner⁴⁵, although its function in B cells has not been
352 known. *RRM2B* may also be associated with repair of DNA damage induced by AID.

353 In summary, our large-scale BCR repertoire analysis provided an overview of B cell abnormalities
354 in IMDs and revealed the presence of new molecules associated with SHM biology in health. This
355 analysis exemplifies the usefulness of multi-omic human data for finding new features in cellular
356 biology.

357



358

359

Fig. 1: Study design and quality features of our dataset.

360

a, Study design used in this report. Sample numbers after quality control are presented under each

361

disease. **b**, Number of RNA-seq reads aligned to BCR regions per million reads in 5 B cell subsets. **c**,

362

Comparison of the per base mutation frequency in the V gene region of 5 B cell subsets. **d**, Comparison

363

of clonotype overlap frequencies between samples according to the individual or the subset being the

364

same or different. “The same individual, the same subset” comparison was performed with SLE

365

samples before and after treatment with belimumab 6 months apart. *P* values were calculated with two-

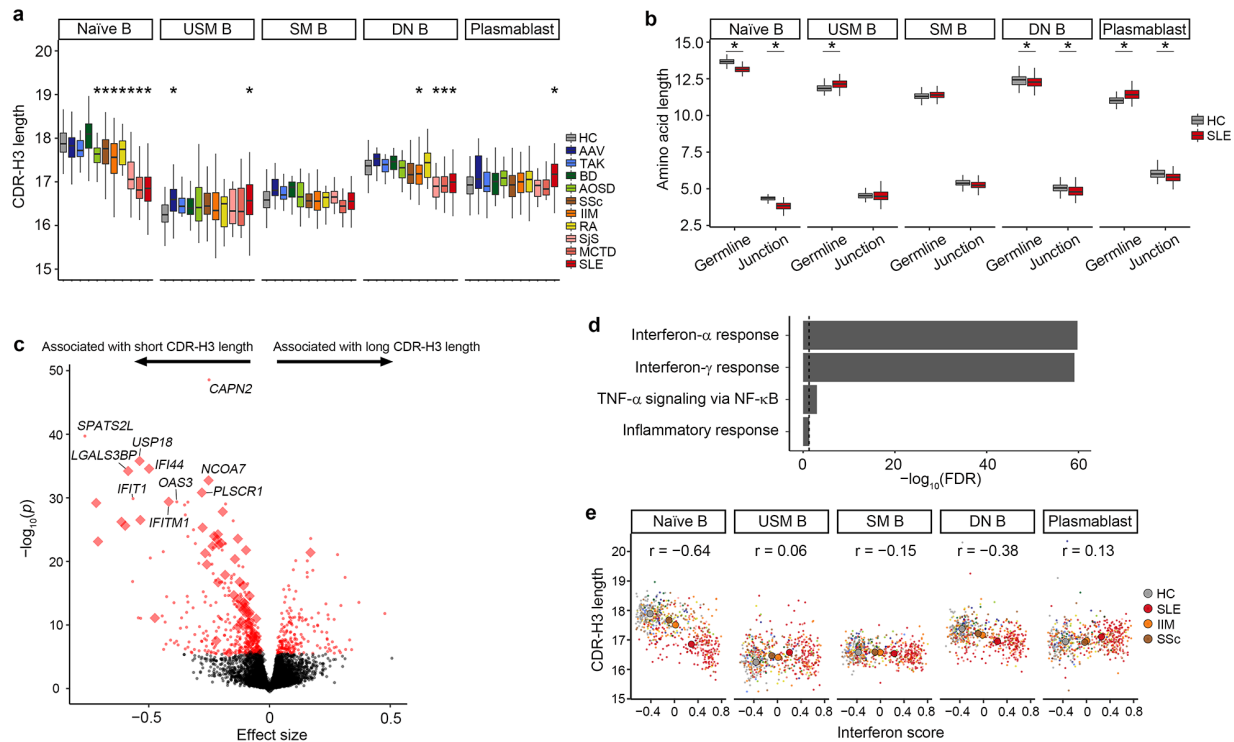
366

sided Mann–Whitney *U*-test. **e**, Comparison of median isotype frequencies in each subset among

367

healthy controls.

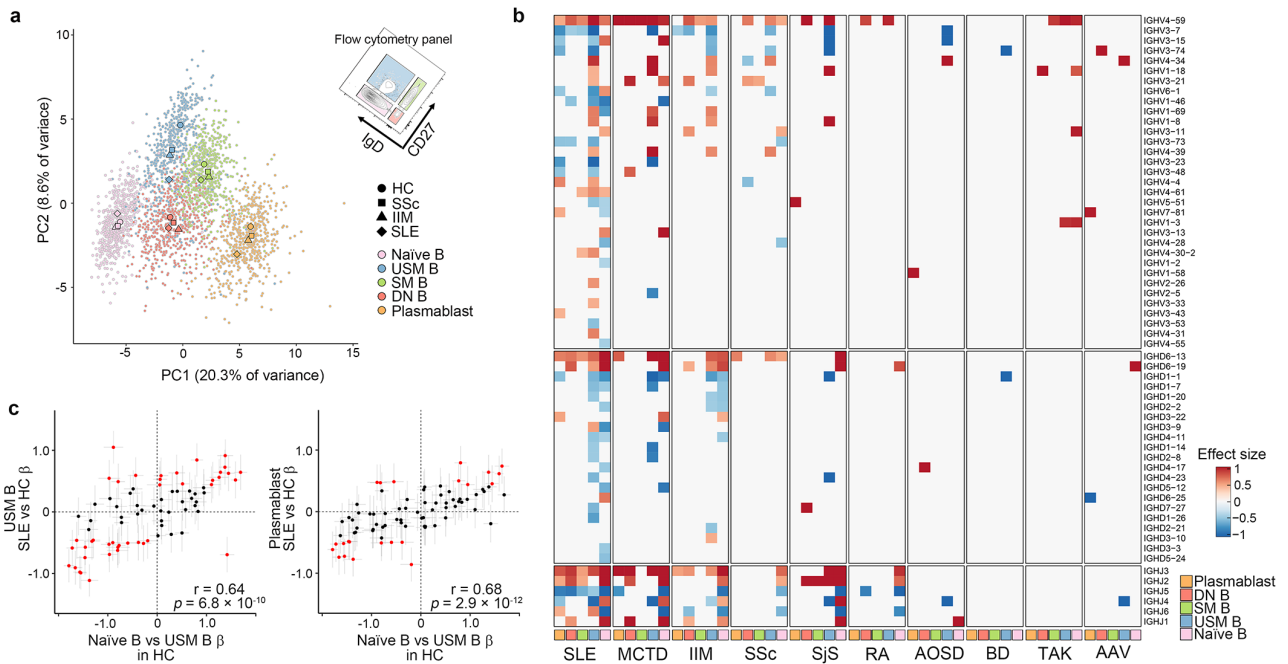
368



369 **Fig. 2: In autoimmune diseases, CDR-H3 length is shortened in naïve B cells in an interferon**
 370 **signal strength-dependent manner.**

371 **a**, Comparison of CDR-H3 length between IMDs and HC. The difference between HC and each disease
 372 was tested with linear regression, adjusting for age and batch effects. *, $p < 0.05/50$. **b**, Comparison of
 373 germline and junction part length of CDR-H3 region (see also Extended data Fig. 2b). The differences
 374 between HC and SLE were tested with linear regression, adjusting for age and batch effects. *, $p <$
 375 $0.05/10$. **c**, Association of naïve B cell gene expression with naïve B cell CDR-H3 length. Fixed effect
 376 meta-analysis p values and effect sizes are plotted (Methods). Genes that attained Bonferroni-adjusted
 377 $p < 0.05$ are marked red. Genes in “HALLMARK_INTERFERON_ALPHA_RESPONSE” gene set
 378 from MSigDB are marked with diamonds. **d**, Gene set enrichment analysis. Results of genes
 379 significantly negatively associated with CDR-H3 length in naïve B cells. FDR, false discovery rate. **e**,
 380 Correlation between IFN score and CDR-H3 length in 5 B cell subsets. Mean values among HC, SLE,
 381 SSc, and IIM are marked with large points. The color of each point indicates the clinical diagnosis as
 382 illustrated in Fig. 1a. Pearson’s r is provided.

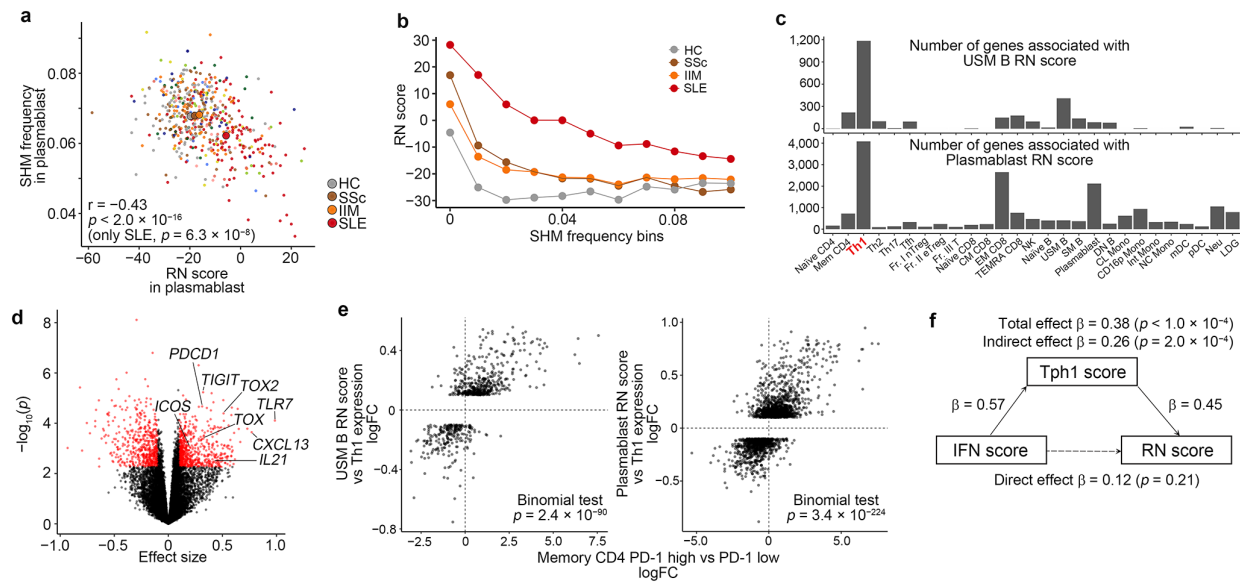
383



384 **Fig. 3: Skewed gene usage in autoimmune diseases suggested an abnormal repertoire maturation**
 385 **process.**

386 **a**, Principal component analysis plot of all samples based on VDJ gene usage. Centroids of HC, SSc,
 387 IIM, and SLE in each subset are marked with large points. **b**, Skewed gene usage in IMDs compared
 388 to HC. In each disease-subset (column), gene usage (row) was compared with HC after adjusting for
 389 age and batch effects. Effect sizes of linear model comparison are plotted if the difference was
 390 statistically significant (FDR < 0.05). Only genes that attained FDR < 0.05 in at least one comparison
 391 are shown. **c**, Comparison of β estimates in usage comparison between naive B vs USM B in HC (y
 392 axis) and SLE vs HC (x axis) in USM B (left) or plasmablasts (right). Each plot corresponds to each
 393 VDJ gene. Genes that attained FDR < 0.05 in SLE vs HC comparison are marked red. Grey bars
 394 indicate 95%CI. Pearson's r and its significance are provided.

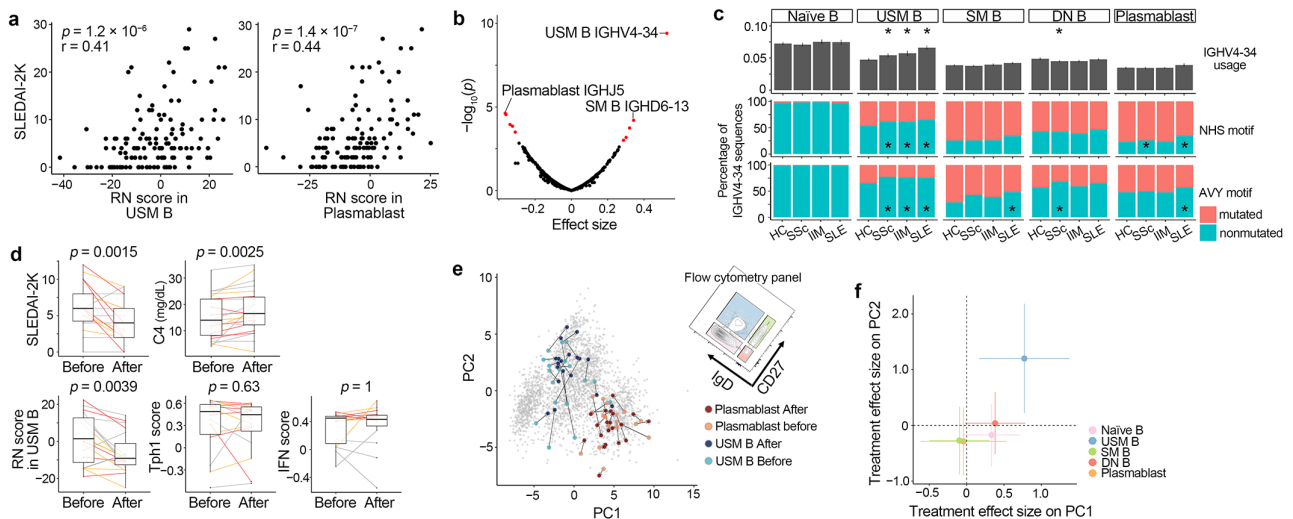
395



396 **Fig. 4: Association of the repertoire abnormality in autoimmune diseases with a low load of**
 397 **mutations and specific subpopulations in Th1 cells.**

398 **a**, Correlation of RN score and SHM frequency in the V gene region of plasmablasts. *P* values were
 399 calculated with linear regression adjusting for diagnosis, age and batch effects. Mean values among
 400 HC, SLE, SSc, and IIM are marked with large points. The color of each point indicates the clinical
 401 diagnosis as illustrated in Fig. 1a. **b**, RN scores of clonotypes stratified by SHM frequency. In each
 402 disease, clonotypes for which ≥ 200 nucleotides in the V gene were sequenced were aggregated and
 403 divided into bins according to mutation frequency in the V gene region. In each bin, RN score was
 404 calculated based on VDJ gene usage. **c**, Number of significantly associated genes with RN score in
 405 USM B (top) and plasmablasts (bottom) among expressed genes in each immune cell. Full description
 406 of subset names is in Table S3. **d**, Association of gene expression in Th1 cells and RN score in USM
 407 B. Significantly associated genes are colored red. **e**, Comparison of the effect size estimates of Th1
 408 expression versus RN score comparison in our cohort (Y-axis) and PD-1^{hi} versus PD-1^{low} T cells
 409 comparison in the data by Rao, *et al*²² (X-axis). Only genes significantly associated with RN scores in
 410 our cohort were included for the analysis. *P* values were calculated by the binomial test for
 411 concordance of the direction of effect size estimates. logFC, log fold change. **f**, Mediation model
 412 representing the relationships among IFN score, Tph1 score and RN score in plasmablasts among SLE
 413 patients. The bootstrapped indirect effect was 0.26 (95% CI, 0.14-0.39; $p = 2 \times 10^{-4}$). Thus, the indirect
 414 effect was statistically significant.

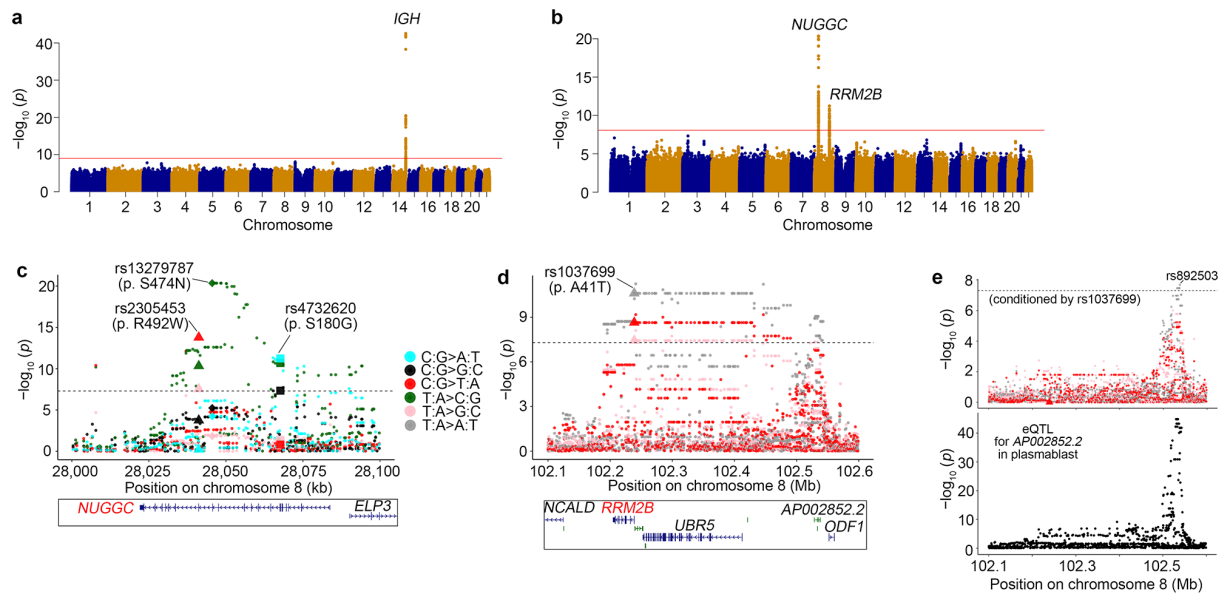
415



416 **Fig. 5: Repertoire abnormalities are associated with disease activity in SLE and attenuated by**
 417 **belimumab treatment.**

418 **a**, Correlation between RN scores and SLEDAI-2K among SLE patients. Pearson's r and its
 419 significance are provided. **b**, Association of VDJ gene usages among 5 B cell subsets and SLEDAI-
 420 2K. Significantly associated genes ($FDR < 0.05$) are red. **c**, Comparison of mean IGHV4-34 usage
 421 (top) and percentage of clonotypes with mutations in the NHS motif (middle) or the AVY motif
 422 (bottom) among IGHV4-34 clonotypes. For usage comparison, the differences between HC and each
 423 disease were tested with linear regression adjusting for age and batch effects. Bars indicate 2SE. For
 424 mutation comparison, Fisher's exact test comparing the number of mutated/nonmutated clonotypes
 425 between HC and each disease was performed in each subset. *, $p < 0.05/15$. **d**, Comparison of clinical
 426 parameters and scores before and after treatment with belimumab. Patients are colored according to
 427 improvement in SLEDAI-2K after treatment (red, ≥ 4 ; orange, 1-3; grey, ≤ 0). Paired Wilcoxon signed-
 428 rank test. **e**, Change in BCR repertoire after belimumab treatment. Paired samples before and after
 429 treatment were connected on the same PCA space as Fig. 3a. **f**, Belimumab treatment effects on PC
 430 scores by fitting the linear mixed model. Bars indicate 95% CI.

431



432 **Fig. 6: Analysis of genetic association reveals molecules associated with somatic hyper-mutation**
433 **biology.**

434 **a**, P values of V gene usage-QTL analysis in naïve B cells. P values of QTL analysis for 43 V genes
435 are plotted together. Red line indicates the significance threshold ($p < 5.0 \times 10^{-8}/43$). **b**, P values of
436 SHM-QTL analysis. P values for QTL analysis for 6 patterns of mutations are plotted together. Red
437 line indicates the significance threshold ($p < 5.0 \times 10^{-8}/6$). **c**, Regional association plot of *NUGGC*
438 locus. Associations of 5 mutation patterns that attained genome-wide significance in this region are
439 plotted. Large marks correspond to 3 independent missense variants that are also (or in tight LD with)
440 lead variants. **d**, Regional association plot of *RRM2B* locus. Associations of 3 mutation patterns that
441 attained genome-wide significance in this region are plotted. A large mark corresponds to the missense
442 variant for *RRM2B* that is in high LD with the lead variant. **e**, Residual association in *RRM2B* locus
443 after conditioning by rs1037699 (top) and eQTL p values for *AP002852.2* expression in plasmablasts
444 (bottom).
445

446 **Methods**

447 **Criteria for patient inclusion/exclusion**

448 The SLE cohort was comprised of patients who had attended the Department of Allergy and
449 Rheumatology at the University of Tokyo Hospital, Division of Rheumatic Diseases at National Center
450 for Global Health and Medicine or Immuno-Rheumatology Center at St. Luke's International Hospital,
451 who met the 1997 revised version of ACR SLE criteria⁴⁶. Of the enrolled 136 patients, 88 cases had
452 active disease with SLEDAI-2K \geq 4. The other cases were inactive at the time of blood withdrawal.
453 Patients administered cyclophosphamide or rituximab within one year or taking \geq 21 mg/day
454 prednisolone were excluded from this analysis. For 22 patients, belimumab was initiated after the
455 collection of the samples, and samples after a half-year treatment were also collected. To assess the
456 drug's effect, cases with concomitant prednisolone increase \geq 5 mg/day or concomitant addition of
457 other immunosuppressants were excluded.

458 The SSc cohort was comprised of patients who had attended the Department of Allergy and
459 Rheumatology at the University of Tokyo Hospital or Department of Rheumatology at Tokyo
460 Metropolitan Komagome Hospital, who met the 2013 American College of Rheumatology (ACR)/
461 European League Against Rheumatism (EULAR) classification criteria for Systemic Sclerosis⁴⁷.
462 Patients taking \geq 21 mg/day of prednisolone were excluded from this analysis.

463 The IIM cohort was comprised of patients who had attended to the Department of Allergy and
464 Rheumatology at the University of Tokyo Hospital or Division of Rheumatology at the Jikei University
465 Hospital. They met standards defined by the Bohan and Peter criteria^{48, 49}, or the European
466 Neuromuscular Center criteria⁵⁰, or the Sontheimer criteria⁵¹ or the Griggs criteria⁵². Patients taking \geq
467 12 mg/day prednisolone were excluded. Of the 85 patients, 36 cases had active disease and required
468 subsequent initiation or an increase of immunomodulatory drugs.

469 The inclusion criteria for HC were people with no apparent co-morbidities, no direct family history

470 of autoimmune diseases and no use of prescription drugs or supplements. Age and sex were matched
471 with the patient cohort as much as possible. For 30 HC cases, only naïve B cells among 5 B cell subsets
472 were collected and utilized for the analysis. For the remaining 107 cases, 5 B cell subsets were
473 collected.

474 Patient inclusion/exclusion criteria for the other diseases were described in detail in our first report².

475 Samples were collected from the University of Tokyo, the Jikei University School of Medicine, St.
476 Luke's International Hospital, National Center for Global Health, and Medicine and Tokyo
477 Metropolitan Komagome Hospital with approval by the Ethics Committees at each site.

478

479 **RNA-seq experiments**

480 We followed the same protocol that we reported previously². Briefly, peripheral blood mononuclear
481 cells (PBMCs) were isolated by density gradient separation with Ficoll-Paque (GE Healthcare)
482 immediately after blood drawing. After lysis of erythrocytes, up to 26 immune cell subsets were
483 collected with purity > 99% using a MoFlo XDP instrument (Beckman Coulter) or a 14-colour cell
484 sorter BD FACSAria Fusion (BD Biosciences). We intended to collect a minimum of 1,000 cells (up
485 to 5,000 cells) per subset. Libraries for RNA-seq were prepared using SMART-seq v4 Ultra Low Input
486 RNA Kits (Takara Bio). Samples failing any of the quality control steps (RNA quality and quantity,
487 PCR amplification, library fragment size) were eliminated from further downstream steps. Libraries
488 were sequenced on the HiSeq 2500 Illumina platform to obtain 100-bp paired-end reads with HiSeq
489 SBS Kit v4 (Illumina).

490

491 **Mapping of the transcriptome**

492 We followed the same protocol we reported previously². Briefly, adaptor sequences were trimmed
493 using cutadapt (v1.16) and 3'- ends with low-quality bases (Phred quality score < 20) were trimmed

494 using the fastx-toolkit (v0.0.14). Reads containing more than 20% low-quality bases were removed.
495 Subsequently, reads were aligned against the GRCh38 reference sequence using STAR (v2.5.3)⁵³ in
496 two-pass mode with GENCODE, version 27 annotations⁵⁴. A median of 10.1 million paired fragments
497 were uniquely mapped per sample and utilized for the subsequent analysis. Expression was quantified
498 using HTSeq (v 0.11.2)⁵⁵.

499

500 **Mapping and QC of BCR sequence**

501 We utilized MiXCR software (version 3.0.13)⁵⁶ for BCR mapping. With FASTQ files after adapter and
502 quality trimming, we utilized the “analyze shotgun” option with parameters “-s hsa”, “--starting-
503 material rna” and “--receptor-type bcr”. After alignment and assembly of clonotypes based on the
504 CDR3 sequence, the full BCR sequence was assembled with the “--contig-assembly” parameter. We
505 set “specificMutationProbability” parameter to 2×10^{-3} . Only productive heavy-chain sequences were
506 kept for the analysis unless specified. For all analyses, samples with more than 500 unique CDR-H3
507 sequences were used for the analysis. For each sample, we treated each unique BCR sequence as a
508 single event, regardless of read depth. For mapping non-productive sequences, we raised the
509 “specificMutationProbability” parameter to 8×10^{-3} in order to reduce the amount of productive
510 sequences with error or mutations counted as non-productive sequences.

511

512 **Mapping of whole genomic sequences**

513 We followed the same protocol that we reported previously². Genomic DNA was isolated from
514 peripheral blood using QIAmp DNA Blood Midi kit (Qiagen) and libraries for whole genome
515 sequencing were prepared using TruSeq DNA PCR-Free Library prep kit (Illumina). Whole genomes
516 were sequenced on Illumina’s HiSeq X with 151-bp pair-end reads. Data processing was based on the
517 standardized best-practice method proposed by GATK (v 4.0.6.0)⁵⁷ with additional recalculation of

518 genotype quality (GQ) scores based on allele frequencies in the 1000 Genomes Project, phase 3, East
519 Asian (EAS) population using CalculateGenotypePosteriors module from GATK. Samples with
520 genotyping call rates < 99% were removed and multi-allelic sites were split into bi-allelic. Calls with
521 GQ < 20 and/or DP < 5 were assigned to missing. Variants with genotyping call rates < 85% and/or
522 HWE P-value < 1.0×10^{-6} were removed. Finally, we used BEAGLE (v 5.1)⁵⁸ to impute missing
523 genotypes. Variants with minor allele frequency < 1% and variants on actionable 59 genes by ACMG⁵⁹
524 were excluded from our analysis because they were out of the scope of this research.

525

526 **Principal component analysis based on VDJ gene usage**

527 All B cell samples used in this study except those samples after belimumab treatment were utilized for
528 PCA analysis. For each sample, each VDJ gene usage was calculated by dividing the number of
529 clonotypes using that gene by the total number of clonotypes in that sample. We retained genes that
530 were used more than 0.5% by 10% or more of the samples. Each gene usage was scaled to mean = 0
531 and variance = 1, and PCA was performed with prcomp function in R.

532 To project samples after belimumab treatment to the same PCA space described above, for each
533 VDJ gene, we centered the mean and standardized the variance with regards to those calculated among
534 all the other B cell samples above. Then, this normalized VDJ gene usage matrix and the PC loadings
535 of each gene calculated in the above analysis were utilized for calculating PC scores for samples after
536 belimumab treatment.

537 We assessed the disease effects on PC scores by linear regression analysis comparing each disease
538 sample versus HC samples in each subset. Treatment effects on PC scores were assessed by the linear
539 mixed model comparing samples before and after treatment, with individuals treated as random effects.
540 The linear mixed model was fit with lme4 package.

541

542 **Calculation of the Repertoire Naïveness Score**

543 For scoring, we only included VDJ genes that were used by more than an average of 0.1% of the
544 clonotypes in each of the 5 subsets. At first, each VDJ gene usage was scaled to mean = 0 and variance
545 = 1 among the 5 B cell subsets of HCs. Then, the effect sizes of VDJ genes were estimated by
546 comparing their usage between naïve B and USM B among HC samples with a linear mixed model,
547 with individuals treated as random effects. Only genes that attained Bonferroni-adjusted p values <
548 0.05 in this comparison were utilized for the subsequent scoring. In addition, genes whose usage
549 showed substantial correlation among 5 B cell subsets of HCs ($|\beta| > 0.4$ in the linear mixed model)
550 were pruned to keep the most significantly different ones in the comparison of USM B and naïve B.
551 As a result, 55 genes (35 V genes, 16 D genes and 4 J genes) were retained for the scoring
552 (Supplementary table 1).

553 For each sample, VDJ gene usage was mean-centered and variance-standardized with regards to
554 those calculated among the 5 B cell subsets of HCs. Then, the following formula was applied for
555 calculating the RN score:

$$556 \quad RNS_i = \sum_{k=1}^{35} \widehat{\beta}_{V_k} \times U_{V_{ik}} + \sum_{k=1}^{16} \widehat{\beta}_{D_k} \times U_{D_{ik}} + \sum_{k=1}^4 \widehat{\beta}_{J_k} \times U_{J_{ik}}$$

557 where $\widehat{\beta}_{V_k}$ corresponded to the effect size of the V_k gene usage comparison between USM B and naïve
558 B among HCs, and $U_{V_{ik}}$ corresponded to the normalized usage of the V_k gene in individual i .

559

560 **Calculation of gene set enrichment score**

561 We utilized the GSVA package⁶⁰ to calculate the gene set enrichment score for samples. After TMM
562 normalization⁶¹ of count data, conversion to log count per million values and batch effect removal with
563 the ComBat function from the sva package⁶², the “gsva” function was utilized for calculating the gene
564 set enrichment score for the gene sets in Supplementary table 2. The Tph1 score was calculated using
565 the 54 genes that were significantly upregulated in PD-1^{hi} CD4⁺ T cells compared with PD-1⁻

566 populations in the original article²² and normalized expression values of Th1 cells. The IFN score was
567 calculated using 97 genes in HALLMARK_INTERFERON_ALPHA_RESPONSE gene set in
568 MSigDB⁶³.

569

570 **Association analysis with gene expression**

571 The association of naïve B cell CDR-H3 length with gene expression in naïve B cells was assessed in
572 each disease at first for reducing the confounding of inter-disease differences in gene expression-CDR-
573 H3 length association. For this step, gene expression data after filtering of low expressed genes was
574 assessed for the association with CDR-H3 length using the limma⁶⁴ package with voom transformation.
575 Known batch effects and age were treated as covariates. Then, resultant effect sizes of each comparison
576 and their standard deviations were utilized for meta-analysis with Metasoft (version 2.0.1) software⁶⁵.

577 For the assessment of association of RN score and gene expression in each immune cell subset, we
578 also utilized limma with voom transformation. The RN score was scaled to mean = 0 and variance =
579 1, and known batch effects and age were treated as covariates. In this comparison, genes fulfilling FDR
580 < 0.05 and $|\log_{2}FC| > 0.1$ were treated as significantly associated genes.

581

582 **Gene set enrichment analysis**

583 We tested the enrichment of CDR-H3 length-associated genes that passed a Bonferroni significance
584 threshold in fixed-effect meta-analysis to hallmark gene sets by MSigDB⁶³ using the clusterProfiler
585 package⁶⁶.

586

587 **Differential usage analysis**

588 For the assessment of disease effects on VDJ gene usage in each subset, we kept genes that were used
589 more than 0.5% by 10% or more of the samples. For each disease, gene usage of the disease and HC

590 samples were scaled to mean = 0 and variance = 1 and tested for the disease effect on gene usage by
591 linear regression with age and known batch effects as covariates.

592

593 **Mediation analysis**

594 The relationships among IFN scores, Tph1 scores and RN scores in plasmablasts were tested among
595 SLE patients by utilizing the classical mediation analysis framework⁶⁷. The IFN score was treated as
596 the independent variable, the Tph1 score was treated as the mediator and the RN score was treated as
597 the outcome variable. P values were calculated via 10000-time bootstrapping using mediation package
598 in R.

599

600 **Calculation of V gene mutation frequency**

601 For the assessment of gene mutation frequency, we quantified per-base mutation frequency in the
602 heavy-chain V gene except for the CDR3 region, that is, in FR1, CDR1, FR2, CDR2 and FR3 regions.
603 As MiXCR assembles the sequence of these regions after alignment and assembly of clonotypes based
604 on CDR3 sequence, only a portion of the nucleotides is actually sequenced in these regions. We
605 counted the number of mutations in these regions in each individual and divided it by the total number
606 of sequenced nucleotides in these regions in the individual to infer the V gene mutation frequency.

607

608 **Comparison with the data from the other cohort**

609 The BCR repertoire sequencing data of peripheral blood mononuclear cells from healthy individual¹¹
610 was downloaded from the NCBI Sequencing Read Archive under BioProject number PRJNA527941.

611 The BCR repertoire sequencing data of sorted B cells from human tonsil¹⁹ was downloaded from
612 ArrayExpress with the accession number E-MTAB-8999.

613 Gene expression data of PD-1^{hi} and PD-1⁻ T cells²² were downloaded from the ImmPort repository

614 with the accession number SDY939. We performed the differential expression analysis between
615 (memoryCD4⁺PD-1^{hi}CXCR5⁻MHCII⁺ICOS⁺ + memoryCD4⁺PD-1^{hi}CXCR5⁻MHCII⁺ICOS⁺) versus
616 memoryCD4⁺PD-1⁻CXCR5⁻ populations to obtain “PD-1^{hi} versus PD-1⁻ logFC”.

617 To validate the reproducibility of the correlation between Tph1 score and RN score in plasmablasts,
618 we utilized peripheral blood Th1 cells and plasmablast RNA-seq data of Japanese SLE patients and
619 HCs that we had obtained previously²⁵, which have no sample overlap with this study.

620

621 **Comparison of the frequency of IGHV4-34 clonotypes with nonmutated AVY/NHS motifs**

622 We counted the number of IGHV4-34 clonotypes with or without mutations in the AVY motif in the
623 FWR1 region or the NHS motif in the CDR2 region. Only clonotypes sequenced at these positions
624 were included in the analysis. As the number of IGHV4-34 clonotypes was not large in one sample,
625 we aggregated the IGHV4-34⁺ clonotypes in each comparison group (e.g., SLE plasmablasts, HC
626 plasmablasts, etc.), and compared the cumulative number of mutated and nonmutated clonotypes
627 between groups by Fisher’s exact test.

628

629 **SHM-QTL analysis**

630 In each individual, all substitutions in V genes in plasmablasts were collected from “mutationsDetailed”
631 columns for FR1, CDR1, FR2, CDR2, and FR3 regions from MiXCR output. The frequency of each
632 of 6 nucleotide substitution pattern (i.e., C:G>A:T, C:G>T:A, C:G>G:C, T:A>C:G, T:A>G:C and
633 T:A>A:T) was calculated in each individual. Nucleotide substitution pattern frequencies were
634 normalized across samples using a rank based inverse normal transformation with the GenABEL
635 package⁶⁸.

636 Association tests of these normalized substitution pattern frequencies and genetic variants were
637 performed using QTLtools (v1.3.1)⁶⁹ with known batch effects, age, sex, disease and the top 2 genetic

638 principal components included as covariates. In total, 509 cases with whole genome information and
639 plasmablast repertoire information were included in this analysis.

640 Variance of substitution patterns explained by 5 lead SHM-QTL variants was estimated by
641 MANOVA analysis using the MVLM package. We calculated the variance of (6 -1) substitution
642 patterns explained by full model (including the 5 lead variants, known batch effects, age and disease
643 term) and a null model (including known batch effects, age and disease term) separately, and the
644 difference between the two values was defined as the variance explained by 5 variants.

645

646 **Data visualization and statistical test**

647 The statistical test performed are indicated in the figure legends or Methods. Throughout the analyses,
648 multiple test correction was performed with Benjamini-Hochberg procedure to obtain corrected q-
649 values unless otherwise indicated. In boxplots, the center lines indicate the median, the box limits
650 indicate the interquartile range (IQR), the whiskers reflect the maximum and minimum values no
651 further than $1.5 \times \text{IQR}$ from the hinge and the points indicate outliers.

652

653 **Acknowledgements**

654 The super-computing resource was provided by Human Genome Center, Institute of Medical Sciences,
655 The University of Tokyo (<http://sc.hgc.jp/shirokane.html>). This study was supported by Chugai
656 Pharmaceutical Co., Ltd., Tokyo, Japan, the Center of Innovation Program from Japan Science and
657 Technology Agency (JST) (JPMJCE1304), the Ministry of Education, Culture, Sports, and the Japan
658 Agency for Medical Research and Development (AMED) (JP17ek0109103, 19ek0410047 and
659 JP20ek0410074).

660 **Author contributions**

661 M.O. conducted bioinformatics analysis with the help of Y.N. and K.I. M.N. contributed to the
662 collection and management of clinical information. M.O., M.N., Y.N., S.K., H.H., R.Y., T.I., Y.A., and
663 H.S. managed and contributed to sample collection, cell sorting, RNA sequencing and whole genome
664 sequencing. H.H., Y.A., and T.I. contributed to QC of the RNA-seq data. A.M. performed protein
665 structure prediction. T.O., K.Y. and K.F designed and managed the project. Y.T. contributed to critical

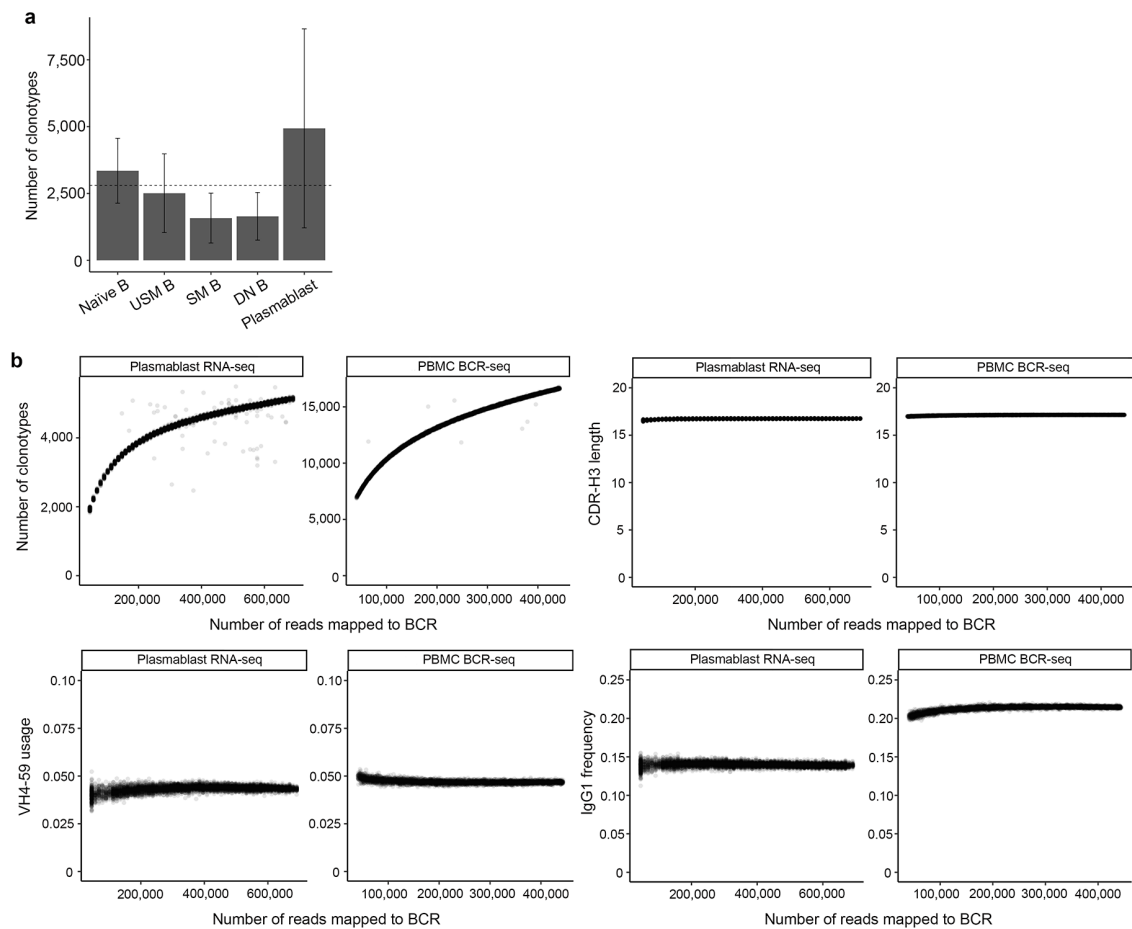
666 reading and revision of the manuscript. M.O. and K.F designed the study and wrote the manuscript
667 with contributions from all authors on the final version of the manuscript.

668 **Competing interests**

669 M.O., Y.N., and T.O. belong to the Social Cooperation Program, Department of functional genomics
670 and immunological diseases, supported by Chugai Pharmaceutical. A.M. is an employee of Chugai
671 Pharmaceutical. K.F. receives consulting honoraria and research support from Chugai Pharmaceutical.

672 **Data availability**

673 BCR data used in this study will be available at the National Bioscience Database Center (NBDC)
674 Human Database (<https://humandbs.biosciencedbc.jp/en/>) at the time of publication. We used publicly
675 available software for the analyses. Custom code is available from the corresponding authors upon
676 reasonable request.
677

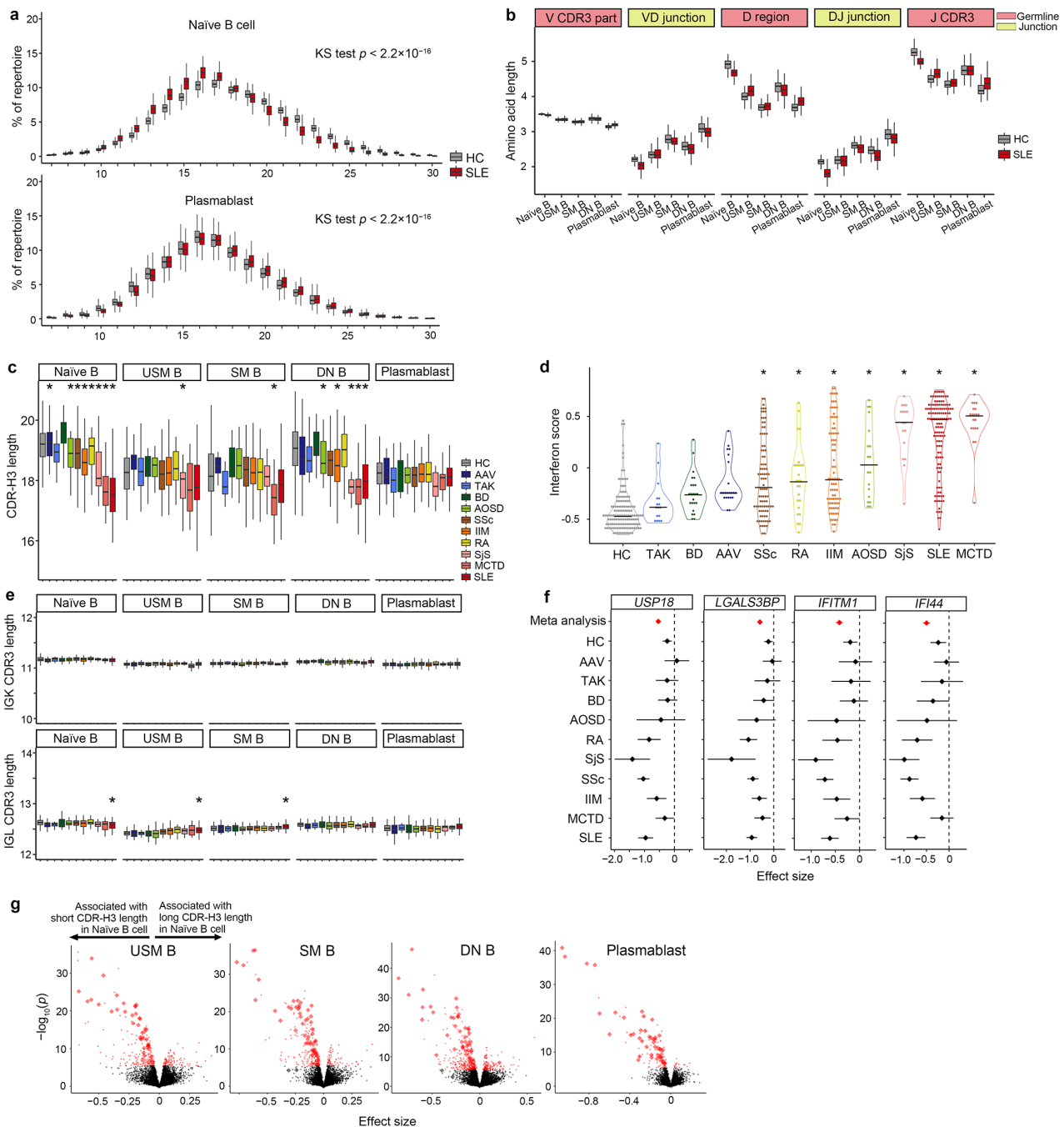


678

679 **Extended data Fig. 1: BCR repertoire characteristics and its association with read depth.**

680 **a**, Mean number of clonotypes identified in each subset in our cohort. Bars indicate 2SD. Dotted line
681 indicates the mean number of identified clonotypes among 5 subsets. **b**, Influence of the number of
682 BCR-mapped reads to the BCR characteristics. The number of reads was down-sampled randomly
683 from the FASTQ files of one plasmablast RNA-seq readout and one peripheral blood mononuclear cell
684 BCR-seq readout¹¹. Changes in the identified number of clonotypes, CDR-H3 length, V gene usage
685 (represented by IGHV 4-59 usage) and isotype frequency (represented by IgG1 frequency) according
686 to the number of reads are plotted.

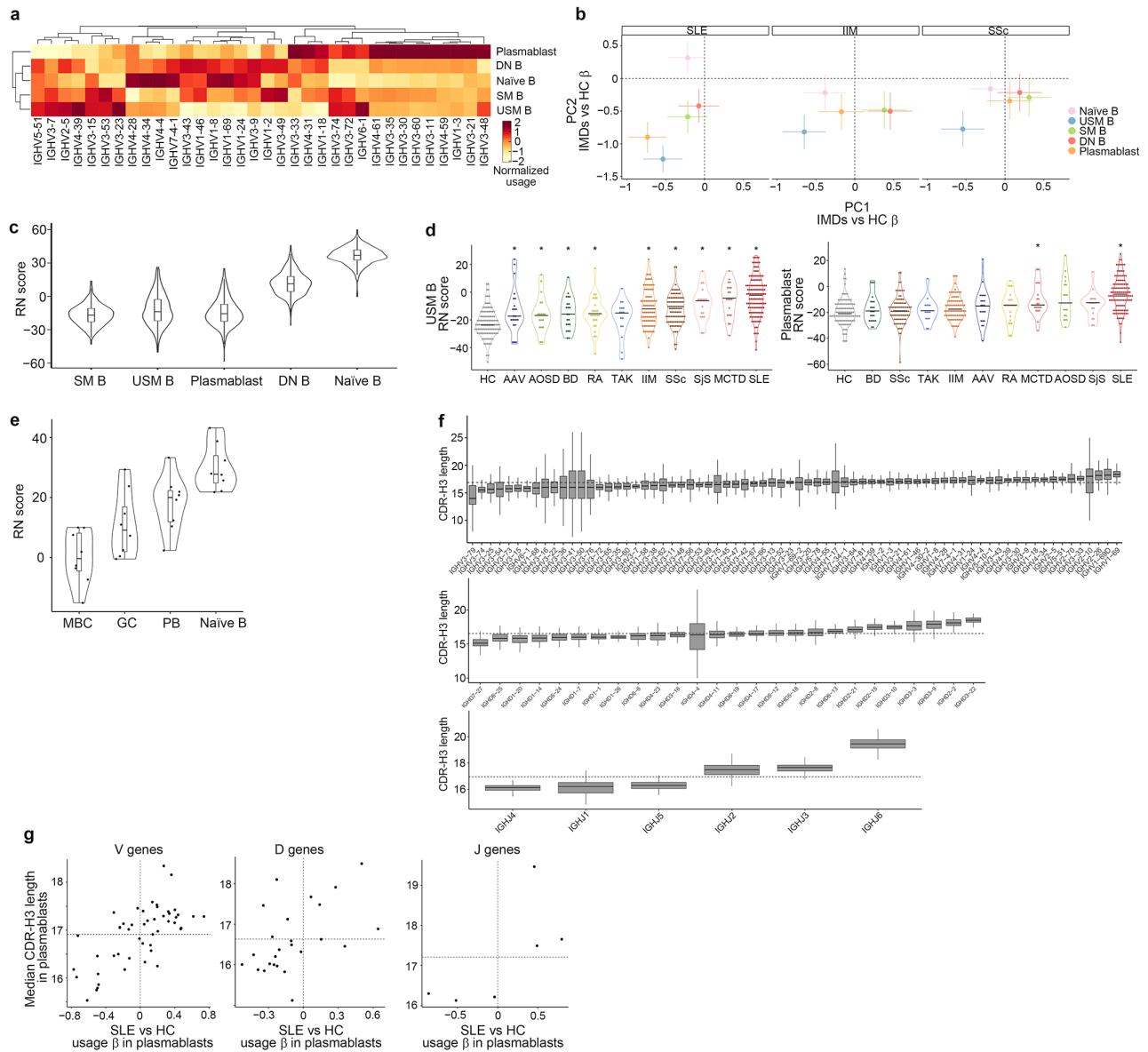
687



688 **Extended data Fig. 2: Association of CDR-H3 length with IMDs and gene expression.**

689 **a**, Comparison of distribution of CDR-H3 lengths in SLE and HC in naïve B cells and plasmablasts. KS test,
 690 Kolmogorov-Smirnov test. **b**, Comparison of lengths of sections that constitute the CDR-H3 region,
 691 related to Fig. 2b. **c**, Comparison of CDR-H3 lengths between IMDs and HC among non-productive
 692 clonotypes. Differences between HC and each disease were tested with linear regression adjusting for
 693 age and batch effects. *, $p < 0.05/50$. **d**, Comparison of IFN scores in naïve B cells among IMD and

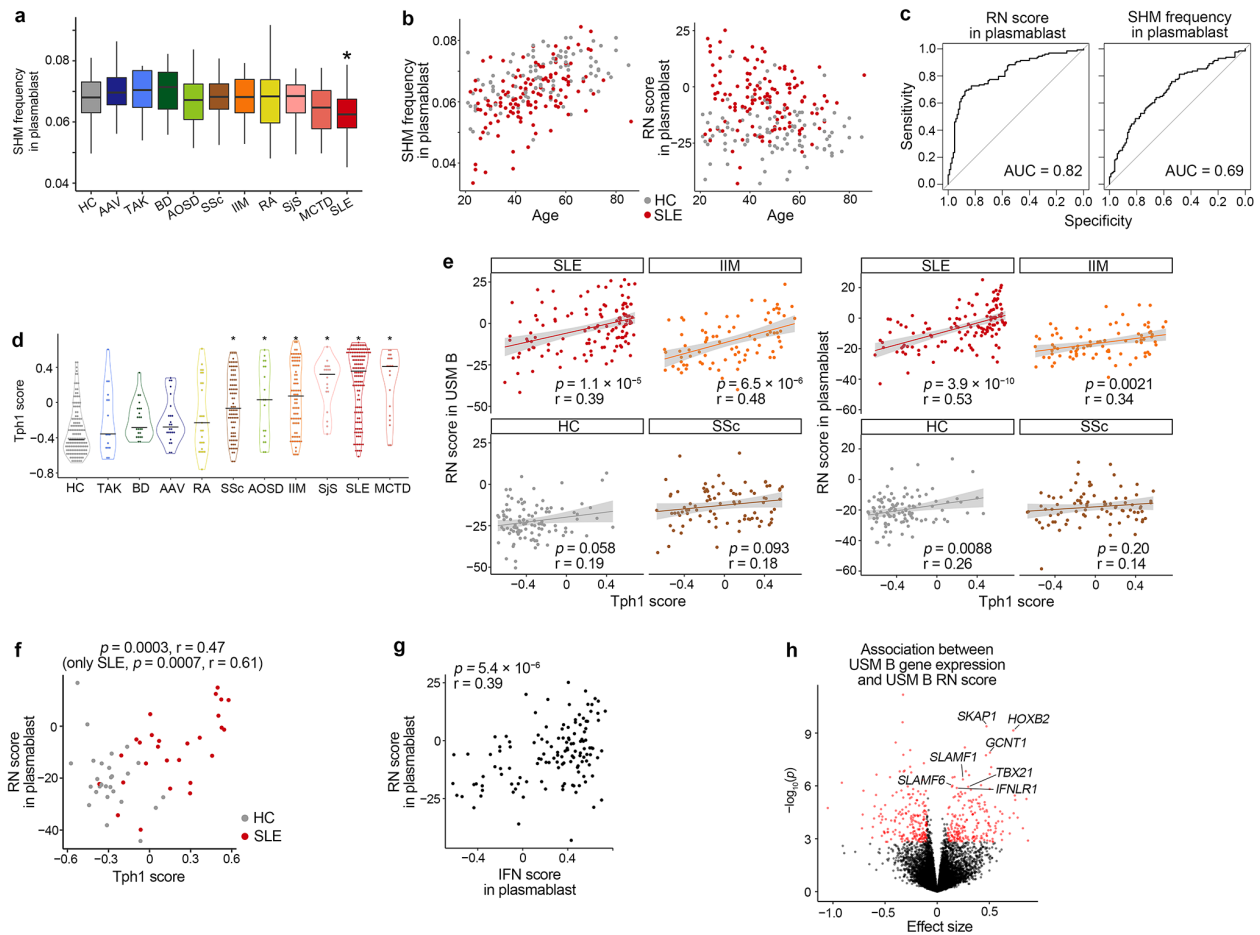
694 HC. Median values are marked with black line. Differences between HC and each disease were tested
695 with linear regression adjusting for age and batch effects. *, $p < 0.05/10$. **e**, Comparison of light chain
696 CDR3 length between IMD and HC. Differences between HC and each disease were tested with linear
697 regression adjusting for age and batch effects. *, $p < 0.05/50$. **f**, Effect sizes for the associations of
698 labeled genes and CDR-H3 lengths in each disease as well as fixed effect meta-analysis. Bars indicate
699 95% CI. **g**, Association of gene expression in 4 B cell subsets with naïve B CDR-H3 length, related to
700 Fig. 2c. Fixed effect meta-analysis p values and effect sizes are plotted. Genes attained Bonferroni-
701 adjusted $p < 0.05$ are marked red. Genes in “HALLMARK_INTERFERON_ALPHA_RESPONSE”
702 gene set are marked as diamonds.
703



704 **Extended data Fig. 3: Comparison of gene usage patterns between IMDs and evaluation of the**
 705 **Repertoire Naïveness score.**

706 **a**, Comparison of mean V gene usage among 5 B cell subsets in HC. Genes that are utilized at least
 707 1% in any of 5 subsets are shown. For comparison, gene usages were Z score normalized column-wise
 708 and hierarchically clustered. **b**, Disease effects on PC scores generated by linear regression analysis:
 709 comparison of each disease sample versus HC samples. Bars indicate 95% CI. **c**, Comparison of RN
 710 scores determined for 5 B cell subsets among all participants in our cohort. **d**, Comparison of RN
 711 scores among IMDs in USM B (left) and plasmablasts (right). Black bars indicate median values.
 712 Differences between HC and each disease were tested with linear regression adjusting for age and
 713 batch effects. *, $p < 0.05/10$. **e**, Comparison of RN scores among previously reported BCR-seq data of

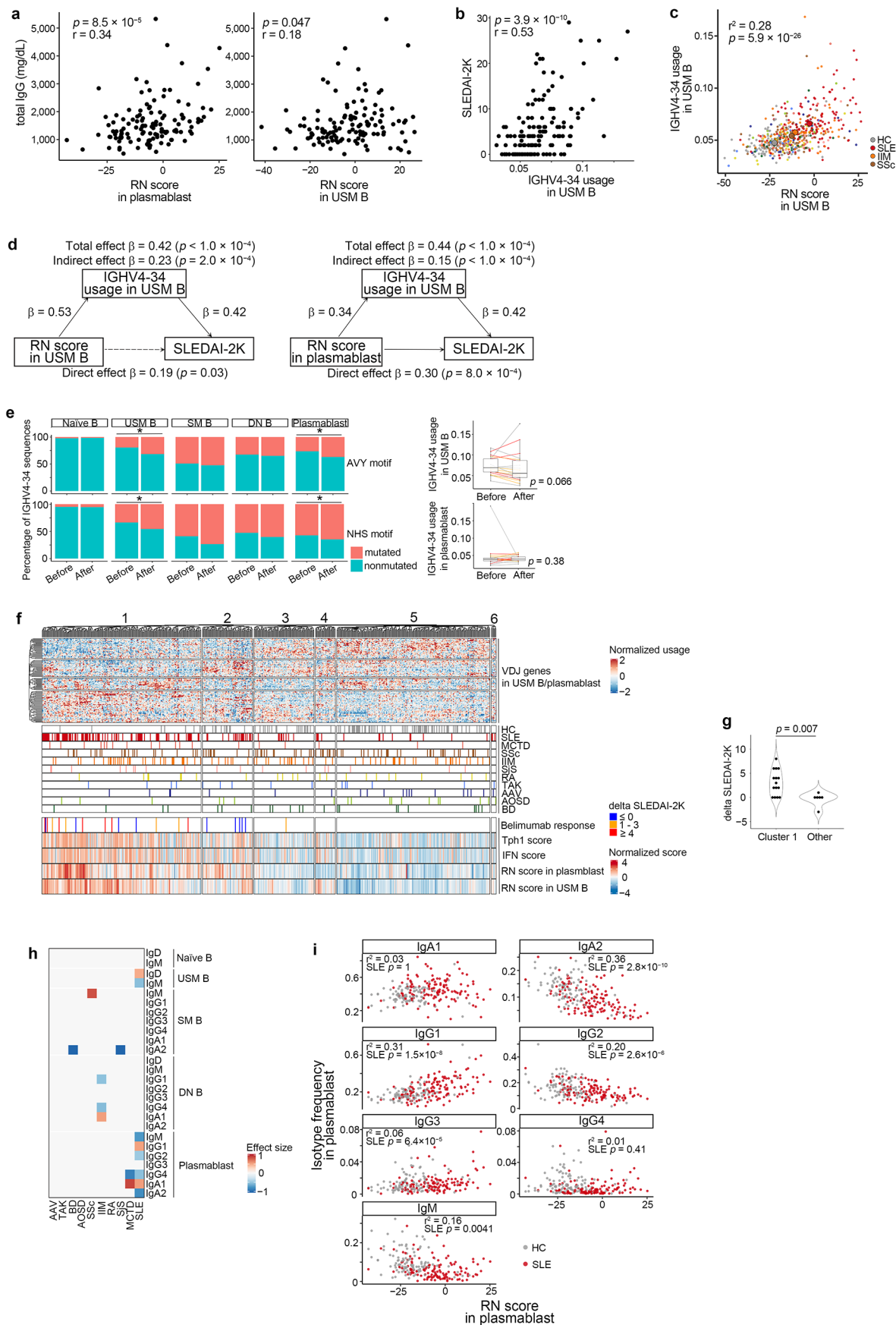
714 sorted B cells from human tonsil. MBC, memory B cell. GC, germinal center B cell. PB, plasmablast.
715 Naïve, naïve B cell. **f**, Distribution of median CDR-H3 lengths of clonotypes utilizing each V (top), D
716 (middle), and J (bottom) gene in plasmablasts among HC. Dotted lines indicate the median CDR-H3
717 length among all clonotypes. **g**, Comparison of the effect size estimates of SLE versus HC gene usage
718 comparison in plasmablasts (x-axis) and median CDR-H3 length of each gene among HC in
719 plasmablasts (y-axis).
720



721 **Extended data Fig. 4: Comparison of RN score-associated gene expression profiles.**

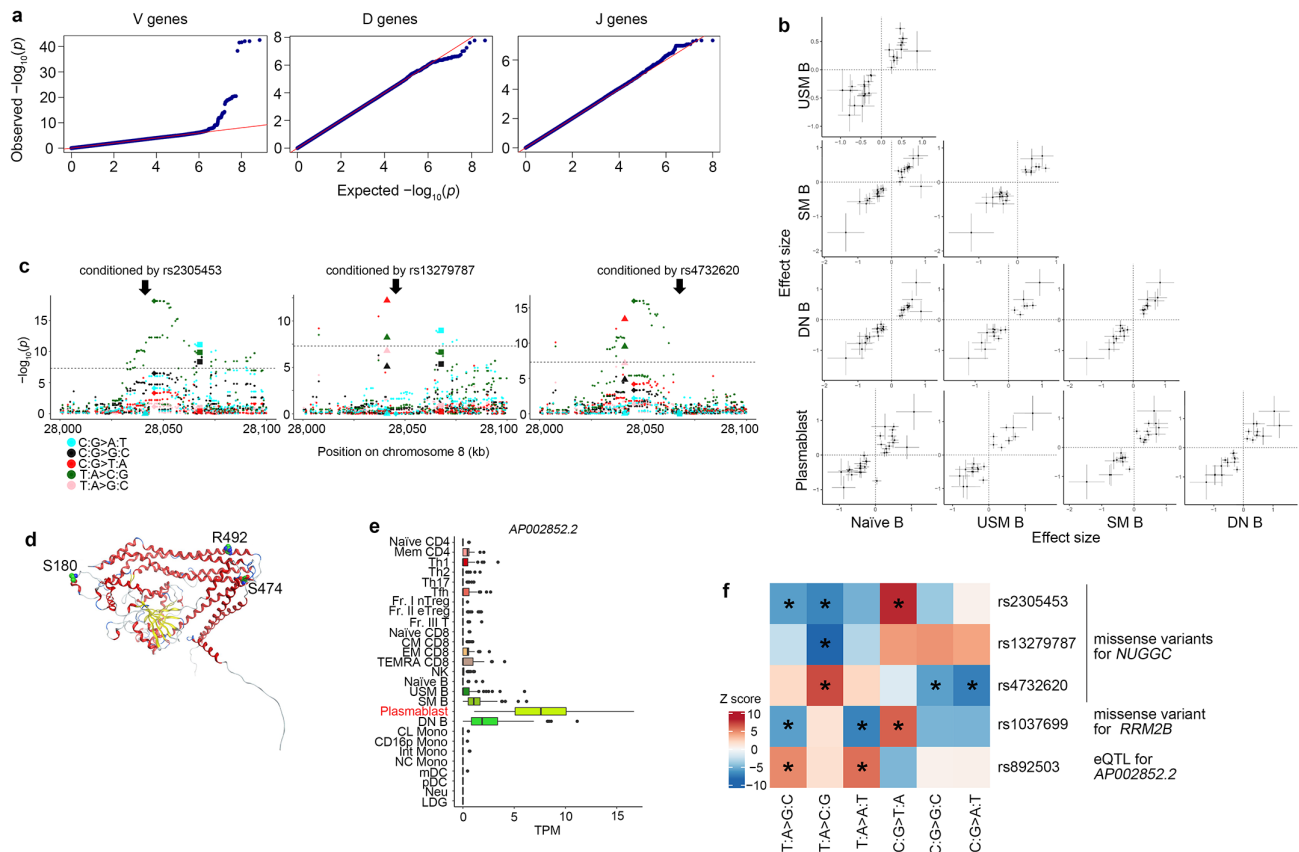
722 **a**, Comparison of SHM frequencies in the V gene region of plasmablasts among IMDs. Differences
 723 between HC and each disease were tested with linear regression adjusting for age and batch effects. *,
 724 $p < 0.05/10$. **b**, Association of SHM frequencies (left) and RN scores of plasmablasts (right) with age.
 725 **c**, ROC curve for classifying SLE patients from HC using RN scores (left) or SHM frequencies (right)
 726 in plasmablasts. **d**, Comparison of Tph1 scores. Differences between HC and each disease were tested
 727 with linear regression adjusting for age and batch effects. *, $p < 0.05/10$. **e**, Correlation between Tph1
 728 scores and RN scores in USM B (left) or plasmablasts (right). Pearson's r and its significance are
 729 provided. **f**, Correlation between Tph1 scores and RN scores in plasmablasts in an independent cohort
 730 of SLE patients and HC. **g**, Correlation between IFN scores and RN scores in plasmablasts among SLE
 731 patients. **h**, Association of gene expression in USM B cells and RN scores in USM B. Significantly
 732 associated genes are marked red.

733



734 Extended data Fig. 5: Comparison of RN score, VH4-34 motifs and isotype frequencies.

735 **a**, Correlation between RN scores and serum total IgG concentrations. Pearson's r and its significance
736 are provided. **b**, Correlation between IGHV4-34 usage in USM B and SLEDAI-2K. Pearson's r and
737 its significance are provided. **c**, Correlation between RN scores and IGHV4-34 usage in USM B. R^2
738 and p values were calculated with linear regression model. Mean values among HC, SLE, SSc, and
739 IIM are marked with large points. The color of each point indicates the clinical diagnosis as illustrated
740 in Fig. 1a. **d**, Mediation model representing the relationships among RN scores in USM B (left) or
741 plasmablasts (right), IGHV4-34 usage in USM B and SLEDAI-2K among SLE patients. **e**, Comparison
742 of the frequencies of clonotypes with mutations in AVY/NHS motifs among IGHV4-34 clonotypes
743 (left) and IGHV4-34 usage (right) before and after treatment with belimumab. For motif comparison,
744 Fisher's exact test comparing the number of mutated/nonmutated clonotypes before and after treatment
745 samples was performed in each subset; *, $p < 0.01$. **f**, Clustering of cases with VDJ gene usage of USM
746 B and plasmablasts that were utilized for RN score calculation. All the participants except for SLE
747 patients after belimumab treatment were hierarchically clustered using average-linkage method based
748 on Pearson's correlation. **g**, Comparison of changes in SLEDAI-2K after belimumab treatment (delta
749 SLEDAI-2K) between SLE patients in cluster 1 and others. **h**, Skewed isotype frequencies in IMDs
750 compared to HC. In each disease (column), isotype frequency (row) was compared with HC after
751 adjusting for age and batch effects. Effect sizes of linear model comparison are plotted if the difference
752 was statistically significant ($FDR < 0.05$). **i**, Correlation between RN scores in plasmablasts and isotype
753 frequency in plasmablast. P values and r^2 values were calculated with linear regression.
754



755 **Extended data Fig. 6: Genetic association analysis with VDJ gene usage and somatic hyper-**
 756 **mutation.**

757 **a**, QQ-plot of V, D, and J gene usage association with genome-wide genetic variants in naïve B cells.
 758 **b**, Comparison of V gene usage-QTL effect sizes among 5 B cell subsets. For each comparison of 2
 759 subsets, the effect sizes were compared for significantly associated V gene-top variant pairs in either
 760 of the subsets. Bars indicate 95%CI. **c**, Regional association plot of the *NUGGC* locus after
 761 conditioning on rs2305453 (left), rs13279787 (middle), and rs4732620 (right). **d**, Estimated protein
 762 structure of NUGGC using the AlphaFold2 algorithm⁷⁰. Three missense variants identified in SHM-
 763 QTL analysis are marked, two of which, Ser180 and R492, are believed to be located on the surface of
 764 the protein. **e**, Expression of *AP002852.2* in immune cells among HCs. **f**, Association of 5 independent
 765 SHM-QTL variants with 6 kinds of mutation patterns. *, $p < 5 \times 10^{-8}$.
 766

767 **Supplementary table legends**

768 **Table S1. Repertoire Naïveness Score calculation.**

769 **Table S2. Gene sets used for calculating gene set enrichment score.**

770 **Table S3. Abbreviations for immune cell subsets.**

771

772 **References**

- 773 1. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317-
774 330 (2015).
- 775
- 776 2. Ota, M. *et al.* Dynamic landscape of immune cell-specific gene regulation in immune-mediated
777 diseases. *Cell* **184**, 3006-3021.e3017 (2021).
- 778
- 779 3. Mitka, M. Treatment for lupus, first in 50 years, offers modest benefits, hope to patients. *Jama*
780 **305**, 1754-1755 (2011).
- 781
- 782 4. Edwards, J.C. *et al.* Efficacy of B-cell-targeted therapy with rituximab in patients with
783 rheumatoid arthritis. *The New England journal of medicine* **350**, 2572-2581 (2004).
- 784
- 785 5. Liston, A., Humblet-Baron, S., Duffy, D. & Goris, A. Human immune diversity: from evolution
786 to modernity. *Nat Immunol* **22**, 1479-1489 (2021).
- 787
- 788 6. van Kempen, T.S., Wenink, M.H., Leijten, E.F., Radstake, T.R. & Boes, M. Perception of self:
789 distinguishing autoimmunity from autoinflammation. *Nat Rev Rheumatol* **11**, 483-492 (2015).
- 790
- 791 7. Elsner, R.A. & Shlomchik, M.J. Germinal Center and Extrafollicular B Cell Responses in
792 Vaccination, Immunity, and Autoimmunity. *Immunity* **53**, 1136-1150 (2020).
- 793
- 794 8. Jenks, S.A. *et al.* Distinct Effector B Cells Induced by Unregulated Toll-like Receptor 7
795 Contribute to Pathogenic Responses in Systemic Lupus Erythematosus. *Immunity* **49**, 725-
796 739.e726 (2018).
- 797
- 798 9. Tipton, C.M. *et al.* Diversity, cellular origin and autoreactivity of antibody-secreting cell
799 population expansions in acute systemic lupus erythematosus. *Nat Immunol* **16**, 755-765
800 (2015).

801

- 802 10. Glass, D.R. *et al.* An Integrated Multi-omic Single-Cell Atlas of Human B Cell Identity.
803 *Immunity* **53**, 217-232.e215 (2020).
804
- 805 11. Ghraichy, M. *et al.* Maturation of the Human Immunoglobulin Heavy Chain Repertoire With
806 Age. *Front Immunol* **11**, 1734 (2020).
807
- 808 12. Miqueu, P. *et al.* Statistical analysis of CDR3 length distributions for the assessment of T and
809 B cell repertoire biases. *Molecular immunology* **44**, 1057-1064 (2007).
810
- 811 13. Gomez-Tourino, I., Kamra, Y., Baptista, R., Lorenc, A. & Peakman, M. T cell receptor β -chains
812 display abnormal shortening and repertoire sharing in type 1 diabetes. *Nat Commun* **8**, 1792
813 (2017).
814
- 815 14. Palanichamy, A. *et al.* Neutrophil-mediated IFN activation in the bone marrow alters B cell
816 development in human and murine systemic lupus erythematosus. *Journal of immunology*
817 (*Baltimore, Md. : 1950*) **192**, 906-918 (2014).
818
- 819 15. Mroczek, E.S. *et al.* Differences in the composition of the human antibody repertoire by B cell
820 subsets in the blood. *Front Immunol* **5**, 96 (2014).
821
- 822 16. Kaplinsky, J. *et al.* Antibody repertoire deep sequencing reveals antigen-independent selection
823 in maturing B cells. *Proceedings of the National Academy of Sciences of the United States of*
824 *America* **111**, E2622-2629 (2014).
825
- 826 17. Bashford-Rogers, R.J.M. *et al.* Analysis of the B cell receptor repertoire in six immune-
827 mediated diseases. *Nature* **574**, 122-126 (2019).
828
- 829 18. Demaison, C., David, D., Fautrel, B. & Theze, J. V(H) gene-family representation in peripheral
830 activated B cells from systemic lupus erythematosus (SLE) patients. *Clinical and experimental*
831 *immunology* **104**, 439-445 (1996).
832
- 833 19. King, H.W. *et al.* Single-cell analysis of human B cell maturation predicts how antibody class
834 switching shapes selection dynamics. *Sci Immunol* **6**, eabe6291 (2021).
835
- 836 20. Sankar, K., Hoi, K.H. & Hötzel, I. Dynamics of heavy chain junctional length biases in
837 antibody repertoires. *Commun Biol* **3**, 207 (2020).

838

- 839 21. William, J., Euler, C., Christensen, S. & Shlomchik, M.J. Evolution of autoantibody responses
840 via somatic hypermutation outside of germinal centers. *Science (New York, N.Y.)* **297**, 2066-
841 2070 (2002).
- 842
- 843 22. Rao, D.A. *et al.* Pathologically expanded peripheral T helper cell subset drives B cells in
844 rheumatoid arthritis. *Nature* **542**, 110-114 (2017).
- 845
- 846 23. Bocharnikov, A.V. *et al.* PD-1hiCXCR5- T peripheral helper cells promote B cell responses in
847 lupus via MAF and IL-21. *JCI Insight* **4**, e130062 (2019).
- 848
- 849 24. Makiyama, A. *et al.* Expanded circulating peripheral helper T cells in systemic lupus
850 erythematosus: association with disease activity and B cell differentiation. *Rheumatology*
851 (*Oxford, England*) **58**, 1861-1869 (2019).
- 852
- 853 25. Takeshima, Y. *et al.* Immune cell multi-omics analysis reveals contribution of oxidative
854 phosphorylation to B cell functions and organ damage of lupus. *bioRxiv*, 2021.10.08.463629
855 (2021).
- 856
- 857 26. Kenderes, K.J. *et al.* T-Bet(+) IgM Memory Cells Generate Multi-lineage Effector B Cells. *Cell*
858 *Rep* **24**, 824-837.e823 (2018).
- 859
- 860 27. Rubtsova, K. *et al.* B cells expressing the transcription factor T-bet drive lupus-like
861 autoimmunity. *J Clin Invest* **127**, 1392-1404 (2017).
- 862
- 863 28. Reed, J.H., Jackson, J., Christ, D. & Goodnow, C.C. Clonal redemption of autoantibodies by
864 somatic hypermutation away from self-reactivity during human immunization. *The Journal of*
865 *experimental medicine* **213**, 1255-1265 (2016).
- 866
- 867 29. Schickel, J.-N. *et al.* Self-reactive VH4-34-expressing IgG B cells recognize commensal
868 bacteria. *The Journal of experimental medicine* **214**, 1991-2003 (2017).
- 869
- 870 30. Banchereau, J. & Pascual, V. Type I interferon in systemic lupus erythematosus and other
871 autoimmune diseases. *Immunity* **25**, 383-392 (2006).
- 872
- 873 31. Tobón, G.J., Izquierdo, J.H. & Cañas, C.A. B lymphocytes: development, tolerance, and their
874 role in autoimmunity-focus on systemic lupus erythematosus. *Autoimmune Dis* **2013**, 827254

- 875 (2013).
876
- 877 32. Rubelt, F. *et al.* Individual heritable differences result in unique cell lymphocyte receptor
878 repertoires of naïve and antigen-experienced cells. *Nat Commun* **7**, 11112 (2016).
879
- 880 33. Schramm, C.A. & Douek, D.C. Beyond Hot Spots: Biases in Antibody Somatic Hypermutation
881 and Implications for Vaccine Design. *Front Immunol* **9**, 1876 (2018).
882
- 883 34. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association
884 studies using summary statistics. *PLoS genetics* **10**, e1004383 (2014).
885
- 886 35. Liu, S. *et al.* Direct measurement of B-cell receptor repertoire's composition and variation in
887 systemic lupus erythematosus. *Genes Immun* **18**, 22-27 (2017).
888
- 889 36. Meffre, E. *et al.* Immunoglobulin heavy chain expression shapes the B cell receptor repertoire
890 in human B cell development. *J Clin Invest* **108**, 879-886 (2001).
891
- 892 37. Wong, E.B., Khan, T.N., Mohan, C. & Rahman, Z.S. The lupus-prone NZM2410/NZW strain-
893 derived Sle1b sublocus alters the germinal center checkpoint in female mice in a B cell-intrinsic
894 manner. *Journal of immunology (Baltimore, Md. : 1950)* **189**, 5667-5681 (2012).
895
- 896 38. Woods, M., Zou, Y.R. & Davidson, A. Defects in Germinal Center Selection in SLE. *Front*
897 *Immunol* **6**, 425 (2015).
898
- 899 39. Yurasov, S. *et al.* Defective B cell tolerance checkpoints in systemic lupus erythematosus. *The*
900 *Journal of experimental medicine* **201**, 703-711 (2005).
901
- 902 40. Cappione, A., 3rd *et al.* Germinal center exclusion of autoreactive B cells is defective in human
903 systemic lupus erythematosus. *J Clin Invest* **115**, 3205-3216 (2005).
904
- 905 41. Muramatsu, M. *et al.* Class switch recombination and hypermutation require activation-
906 induced cytidine deaminase (AID), a potential RNA editing enzyme. *Cell* **102**, 553-563 (2000).
907
- 908 42. Pilzecker, B. & Jacobs, H. Mutating for Good: DNA Damage Responses During Somatic
909 Hypermutation. *Front Immunol* **10**, 438 (2019).
910
- 911 43. Richter, K. *et al.* Speckled-like pattern in the germinal center (SLIP-GC), a nuclear GTPase

- 912 expressed in activation-induced deaminase-expressing lymphomas and germinal center B cells.
913 *J Biol Chem* **284**, 30652-30661 (2009).
914
- 915 44. Richter, K. *et al.* Altered pattern of immunoglobulin hypermutation in mice deficient in Slip-
916 GC protein. *J Biol Chem* **287**, 31856-31865 (2012).
917
- 918 45. Tanaka, H. *et al.* A ribonucleotide reductase gene involved in a p53-dependent cell-cycle
919 checkpoint for DNA damage. *Nature* **404**, 42-49 (2000).
920
- 921 46. Hochberg, M.C. Updating the American College of Rheumatology revised criteria for the
922 classification of systemic lupus erythematosus. *Arthritis Rheum* **40**, 1725 (1997).
923
- 924 47. van den Hoogen, F. *et al.* 2013 classification criteria for systemic sclerosis: an American
925 College of Rheumatology/European League against Rheumatism collaborative initiative.
926 *Arthritis Rheum* **65**, 2737-2747 (2013).
927
- 928 48. Bohan, A. & Peter, J.B. Polymyositis and dermatomyositis (first of two parts). *The New*
929 *England journal of medicine* **292**, 344-347 (1975).
930
- 931 49. Bohan, A. & Peter, J.B. Polymyositis and dermatomyositis (second of two parts). *The New*
932 *England journal of medicine* **292**, 403-407 (1975).
933
- 934 50. Hoogendijk, J.E. *et al.* 119th ENMC international workshop: trial design in adult idiopathic
935 inflammatory myopathies, with the exception of inclusion body myositis, 10-12 October 2003,
936 Naarden, The Netherlands. *Neuromuscul Disord*, vol. 14: England, 2004, pp 337-345.
937
- 938 51. Sontheimer, R.D. Would a new name hasten the acceptance of amyopathic dermatomyositis
939 (dermatomyositis sine myositis) as a distinctive subset within the idiopathic inflammatory
940 dermatomyopathies spectrum of clinical illness? *J Am Acad Dermatol* **46**, 626-636 (2002).
941
- 942 52. Griggs, R.C. *et al.* Inclusion body myositis and myopathies. *Ann Neurol* **38**, 705-713 (1995).
943
- 944 53. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics (Oxford, England)*
945 **29**, 15-21 (2013).
946
- 947 54. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic*
948 *acids research* **47**, D766-d773 (2019).

949

- 950 55. Anders, S., Pyl, P.T. & Huber, W. HTSeq--a Python framework to work with high-throughput
951 sequencing data. *Bioinformatics (Oxford, England)* **31**, 166-169 (2015).
952
- 953 56. Bolotin, D.A. *et al.* MiXCR: software for comprehensive adaptive immunity profiling. *Nat*
954 *Methods* **12**, 380-381 (2015).
955
- 956 57. DePristo, M.A. *et al.* A framework for variation discovery and genotyping using next-
957 generation DNA sequencing data. *Nature genetics* **43**, 491-498 (2011).
958
- 959 58. Browning, B.L., Zhou, Y. & Browning, S.R. A One-Penny Imputed Genome from Next-
960 Generation Reference Panels. *Am J Hum Genet* **103**, 338-348 (2018).
961
- 962 59. Kalia, S.S. *et al.* Recommendations for reporting of secondary findings in clinical exome and
963 genome sequencing, 2016 update (ACMG SF v2.0): a policy statement of the American
964 College of Medical Genetics and Genomics. *Genetics in Medicine* **19**, 249-255 (2017).
965
- 966 60. Hänzelmann, S., Castelo, R. & Guinney, J. GSVA: gene set variation analysis for microarray
967 and RNA-seq data. *BMC Bioinformatics* **14**, 7 (2013).
968
- 969 61. Robinson, M.D., McCarthy, D.J. & Smyth, G.K. edgeR: a Bioconductor package for
970 differential expression analysis of digital gene expression data. *Bioinformatics (Oxford,*
971 *England)* **26**, 139-140 (2010).
972
- 973 62. Leek, J.T., Johnson, W.E., Parker, H.S., Jaffe, A.E. & Storey, J.D. The sva package for
974 removing batch effects and other unwanted variation in high-throughput experiments.
975 *Bioinformatics (Oxford, England)* **28**, 882-883 (2012).
976
- 977 63. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection.
978 *Cell Syst* **1**, 417-425 (2015).
979
- 980 64. Ritchie, M.E. *et al.* limma powers differential expression analyses for RNA-sequencing and
981 microarray studies. *Nucleic acids research* **43**, e47 (2015).
982
- 983 65. Han, B. & Eskin, E. Interpreting meta-analyses of genome-wide association studies. *PLoS*
984 *genetics* **8**, e1002555 (2012).
985

- 986 66. Yu, G., Wang, L.G., Han, Y. & He, Q.Y. clusterProfiler: an R package for comparing biological
987 themes among gene clusters. *Omic*s **16**, 284-287 (2012).
988
- 989 67. Baron, R.M. & Kenny, D.A. The moderator-mediator variable distinction in social
990 psychological research: conceptual, strategic, and statistical considerations. *J Pers Soc Psychol*
991 **51**, 1173-1182 (1986).
992
- 993 68. Aulchenko, Y.S., Ripke, S., Isaacs, A. & van Duijn, C.M. GenABEL: an R library for genome-
994 wide association analysis. *Bioinformatics (Oxford, England)* **23**, 1294-1296 (2007).
995
- 996 69. Delaneau, O. *et al.* A complete tool set for molecular QTL discovery and analysis. *Nat Commun*
997 **8**, 15452 (2017).
998
- 999 70. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583-
1000 589 (2021).
1001
1002