

1 Synthetic standards combined with error and bias correction 2 improves the accuracy and quantitative resolution of antibody 3 repertoire sequencing in human naïve and memory B cells

4
5 Simon Friedensohn*¹, John M. Lindner*², Vanessa Cornacchione², Mariavittoria Iazeolla², Enkelejda Miho¹, Andreas
6 Zingg¹, Simon Meng¹, Elisabetta Traggiai², and Sai T. Reddy¹

7
8 ¹Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland

9 ²Novartis Institutes for BioMedical Research, Basel, Switzerland

10 *equal contribution

11
12 Correspondence: elisabetta.traggiai@novartis.com and sai.reddy@ethz.ch

13 14 ABSTRACT

15
16 **High-throughput sequencing of immunoglobulin repertoires (Ig-seq) is a powerful method for**
17 **quantitatively interrogating B cell receptor sequence diversity. When applied to human**
18 **repertoires, Ig-seq provides insight into fundamental immunological questions, and can be**
19 **implemented in diagnostic and drug discovery projects. However, a major challenge in Ig-seq**
20 **is ensuring accuracy, as library preparation protocols and sequencing platforms can**
21 **introduce substantial errors and bias that compromise immunological interpretation. Here,**
22 **we have established an approach for performing highly accurate human Ig-seq by combining**
23 **synthetic standards with a comprehensive error and bias correction pipeline. First, we**
24 **designed a set of 85 synthetic antibody heavy chain standards (*in vitro* transcribed RNA) to**
25 **assess correction workflow fidelity. Next, we adapted a library preparation protocol that**
26 **incorporates unique molecular identifiers (UIDs) for error and bias correction which, when**
27 **applied to the synthetic standards, resulted in highly accurate data. Finally, we performed Ig-**
28 **seq on purified human circulating B cell subsets (naïve and memory), combined with a**
29 **cellular replicate sampling strategy. This strategy enabled robust and reliable estimation of**
30 **key repertoire features such as clonotype diversity, germline segment and isotype subclass**
31 **usage, and somatic hypermutation (SHM). We anticipate that our standards and error and**
32 **bias correction pipeline will become a valuable tool for researchers to validate and improve**
33 **accuracy in human Ig-seq studies, thus leading to potentially new insights and applications in**
34 **human antibody repertoire profiling.**

35 36 37 INTRODUCTION

38
39 Adaptive immune responses are governed by cooperative interactions between B and T
40 lymphocytes upon antigen recognition. A hallmark of these cells is the somatic generation of
41 clonally unique antigen receptors during primary lymphocyte differentiation. In particular, B cell
42 antigen receptors (BCRs, and their analogous secreted form, antibodies) result from rearrangement
43 of the germline-encoded variable (V), diversity (D, heavy chain only), and joining (J) gene
44 segments. V(D)J recombination in B cells creates a highly complex receptor population (generally
45 interchangeably referred to as BCR, antibody, or immunoglobulin (Ig) repertoires), which matures
46 upon antigen experience to produce the more targeted, high-affinity memory BCR network. In-
47 depth and accurate characterization of these repertoires provides valuable insight into the generation

48 and maintenance of immunocompetency, which can be used to monitor changes in immune status,
49 and to identify potentially reactive clones for therapeutic or other uses. Due to rapid technological
50 advances, high-throughput sequencing of Ig genes (Ig-seq) has become a major approach to catalog
51 the diversity of antibody repertoires (1-3). Ig-seq applied to human B cells has potential in a variety
52 of applications (4), particularly in antibody drug discovery (5-7), diagnostic profiling for vaccines
53 (8, 9), and biomarker-based disease detection (10, 11). Additionally, Ig-seq is enabling a more
54 comprehensive understanding of basic human immunobiology, such as B cell clonal distribution
55 across physiological compartments in health and disease (12, 13).

56
57 A major challenge in Ig-seq is the requirement of accurate and high-quality datasets. Several current
58 library preparation protocols are based on target enrichment from genomic DNA or mRNA (14).
59 For example, the conversion of mRNA (more commonly used due to transcript abundance and
60 isotype splicing) into antibody sequencing libraries relies on a number of molecular reagents and
61 amplification steps (e.g. reverse transcriptase, multiplex primer sets, PCR), which potentially
62 introduce errors and bias. Due to the highly polymorphic nature of repertoires especially from
63 affinity-matured memory B cells and plasma cells, it becomes essential to determine if such
64 technical noise occurs at non-negligible rates, as this could alter quantitation of critical repertoire
65 features such as clonal frequencies, germline gene usage, and somatic hypermutation (SHM) (14,
66 15). One way to address this is through the use of synthetic control standards, for which the
67 sequence and abundance is known prior to sequencing, thus providing a means to assess quality and
68 accuracy (16). Several examples of standards have been presented for Ig-seq; Shugay et al.
69 sequenced libraries prepared from a small polyclonal pool of B and T lymphocyte cell lines, and
70 observed nearly 5% erroneous reads, resulting in approximately 100 false-positive variants per
71 clone (17). Recently, Khan et al. developed a set of synthetic RNA (*in vitro* transcribed) spike-in
72 standards based on mouse antibody sequences, which were used to show that a substantial amount
73 of errors and bias are introduced during multiplex-PCR library preparation and sequencing (18).

74 Various experimental and computational workflows exist to mitigate the effects of errors and bias
75 in Ig-seq. One of the most advanced and powerful strategies is to prepare libraries with the
76 incorporation of random and unique molecular identifiers (UIDs, also commonly referred to as
77 UMIs or molecular barcodes). Following sequencing, error correction can be performed by
78 clustering and consensus building of reads that share the same UID; reads sharing the same UID are
79 assumed to be derived from the same original mRNA/cDNA molecule (19). Furthermore, bias
80 correction for cDNA abundance can be performed by counting the number of UID (instead of total
81 reads) (20, 21). Several iterations of UID-tagging have been developed for Ig-seq, such as UID
82 labeling during first- and second-strand cDNA synthesis (22), UID addition during RT template
83 switching (23), and so-called “tagmentation” of UID-labelled amplicons (24). Recently, we
84 developed an innovative strategy to add UID both during first-strand cDNA synthesis as well as
85 multiplex-PCR amplification; this protocol, known as molecular amplification fingerprinting
86 (MAF), results in comprehensive error and bias correction of mouse antibody repertoires (18).

87 Here, we describe an experimental-computational approach to generate highly accurate human Ig-
88 seq data. We first designed a comprehensive set of synthetic standards based on human antibody
89 heavy chain variable (IGHV) sequences: a total of 85 *in vitro* transcribed RNA standards, each with
90 a unique complementarity determining region 3 (CDR3) sequence and covering nearly the entire set
91 of productive human Ig heavy chain (IgH) germline (IGHV) gene segments. We used these

92 synthetic standards to quantify the impact of errors and bias introduced during multiplex-PCR
93 library preparation, and the robustness with which our previously developed method for UID
94 addition by MAF could correct these artifact sequences. Finally, we implemented MAF-based error
95 and bias correction on human B cell subsets (naïve and memory), which enabled us to make
96 accurate clonal diversity estimates and quantify divergent repertoire features across B cell
97 compartments.

98

99 **RESULTS**

100 **Design of a comprehensive set of human synthetic standards**

101 Our previously established set of murine synthetic antibody standards contained 16 unique clones
102 (CDR3s) covering 7 IGHV gene segments (out of more than 140 annotated murine IGHV gene
103 segments) (18). For our human standards, we developed a more comprehensive set consisting of 85
104 clones encompassing nearly the entire germline IGHV repertoire. The most commonly used
105 repository for human germline segments is the International ImmunoGenetics Database (IMGT),
106 which has annotated 61 IGHV alleles as functional or having an open reading frame (25). After
107 filtering out paralogs and selecting only gene segments that have been found in productive
108 rearrangements (2), we chose 48 IGHV gene segments as the basis for our standards (**Table S1**).
109 Each standard contained the following elements (5' to 3'): (i) a conserved non-coding region, (ii)
110 ATG start codon and a leader peptide sequence spliced to its respective IGHV gene segment, (iii) a
111 synthetic CDR3 sequence, (iv) a germline IgH J (IGHJ) gene segment, (v) a non-coding synthetic
112 sequence identifier (for the separation of standards from biological sequences), (vi) a partial
113 segment of the constant region from isotypes IgM, IgG, and IgA (**Figure 1A**). This design allows
114 amplification of our synthetic controls with a variety of PCR primer sets. Notably, for control
115 singleplex-PCR experiments, all standards can be amplified by a single forward primer (targeting
116 the conserved 5' non-coding region) and a single reverse primer (targeting one of the isotype
117 constant regions). Since IGHV gene segment usage has been reported to be non-uniform (2, 26), we
118 selected the most abundant segments for use in multiple standards (**Figure 1B, Table S1**).

119 All standards carry a unique CDR3 sequence, which visually aids the analysis of sequencing results.
120 Furthermore, all clones were designed to be resilient against sequencing and PCR errors: at least 9
121 specific nucleotide (nt) deletions, insertions, and/or mutations are needed in order to turn one CDR3
122 nt sequence into another (**Figure 1B**). For our experiments, synthetic standard genes were *in vitro*
123 transcribed to RNA and subsequently reverse transcribed to cDNA. We measured individual cDNA
124 molecules by digital droplet PCR (ddPCR) and capillary electrophoresis. Standards were then
125 pooled in a non-uniform concentration distribution and maintained as a master stock (**Table S1**).

126 **Human Ig-seq library preparation using the MAF protocol**

127 We adapted our previously described library preparation protocol for murine antibody repertoires to
128 be compatible for human Ig-seq (18). This protocol is based on targeted amplification via RT of
129 RNA to first-strand cDNA, followed by two PCR amplification steps (27, 28), the first of which
130 uses a forward multiplex primer set targeting the IGHV framework region 1 (FR1). Each step also

131 incorporates fragments of Illumina sequencing adapters (IA), such that the final product of the
132 workflow is already compatible with the Illumina sequencing platform (**Figure 2A**).

133 Importantly, our library preparation protocol used primers incorporating random-nucleotide UIDs,
134 thus enabling MAF-based error and bias correction. A reverse-UID (RID) with theoretical diversity
135 up to 2×10^7 unique sequences is present in the RT primer (between the Ig constant region-specific
136 and partial IA regions), and a forward-UID (FID with additional diversity of approximately 7×10^5
137 unique sequences) is present on the forward multiplex primer set (between the FR1-specific and
138 partial IA regions) used in the first PCR reaction. Such high diversity among RIDs is necessary in
139 order to prevent tagging of different cDNA molecules with the same barcode. Our multiplex
140 forward primer set was designed to target all IMGT-annotated IGHV gene segments (**Figure 2B**).
141 To compromise between maintaining similar amplicon length across gene segment families and
142 creating thermodynamically equivalent oligonucleotides, we placed the primers at or near the
143 beginning of each FR1, resulting in a melting temperature range of 57°C to 63°C (**Figure 2C**). This
144 range and the accompanying (unavoidable) variability in GC content have been shown to
145 potentially cause differences in amplification efficiencies, which in turn leads to a biased
146 representation of segment usage frequencies in Ig-seq data. Our workflow aimed to solve this
147 problem in two ways: first, since the RID labels cDNA at the single molecule level, we are able to
148 resolve the number of molecules by counting the number of RIDs instead of raw reads. Second, by
149 using the FIDs on our forward primers, we are able to further normalize our molecular count, since
150 Ig genes preferentially amplified by our primer set should show a higher ratio of FIDs to RIDs.
151 Additionally, the RIDs can be used to correct for errors introduced during PCR and the sequencing
152 process itself by grouping sequencing reads based on their RID, then correcting diverging nt
153 positions by generating a consensus sequence (majority voting scheme). This is especially useful in
154 Ig-seq when attempting to distinguish true SHM variants from erroneous sequences.

155 **Combining standards with MAF to correct errors and bias in Ig-seq**

156 To evaluate the extent of errors and bias present in human Ig-seq data, standards were mixed
157 (spiked-in) with cDNA prepared from circulating purified human B cells. In total, we sequenced 28
158 independently prepared libraries and annotated them with a custom aligner (18, 29). Prior to
159 alignment, reads were either kept as uncorrected (raw) reads or were corrected using our MAF
160 pipeline that takes into account RID and FID information. In this way, we could directly compare
161 the number of erroneous sequences produced in uncorrected vs. MAF-corrected datasets. Clonal
162 assignment of uncorrected reads produced many erroneous CDR3 amino acid (a.a.) variants
163 (sequences with at least 1 a.a. difference from the nearest standard control sequence); for example,
164 in a dataset with 100,000 aligned reads, there was a median value of 23 errors per clone (**Figure**
165 **3A**). The number of erroneous variants produced showed a clear correlation with the individual
166 abundance of each clone within the master stock ($r = 0.89$). When taking the entire VDJ nt
167 sequence into account, an even greater number of erroneous variants were observed (≥ 1 nt
168 difference from the standard sequence) (**Figure 3B**). We observed a median value of 118 erroneous
169 nt variants per standard (per 100,000 aligned sequences). Again, the number of erroneous variants
170 exhibited a clear correlation with clone abundance ($r = 0.90$). However, we did not observe any
171 significant trend linking IGHV family to the error rate (F-Test on full and reduced linear model, $p =$
172 0.083).

173 After performing error correction with our MAF pipeline, there was a dramatic reduction of CDR3
174 and VDJ errors: we observed a median value of 0 and 1 error per clone, respectively. Across all
175 datasets, we remove an average of 94.2% CDR3 a.a. and 97.4% VDJ nt erroneous variants (**Figure**
176 **3A-B**). For example, prior to error correction the standard “CARGINGERALEW” (from dataset
177 Donor 1, IgG aliquot 1, see **Table S3**) displayed 39 additional CDR3 a.a. variants and 217
178 additional VDJ nt variants (**Figure 3D**). After error correction, we retain only the correct CDR3 a.a.
179 sequence and only one additional (erroneous) nt variant. With our current filtering criteria (see
180 Methods), we observed 16 instances in which we removed a standard control sequence that was
181 present in the raw data. However, in the vast majority of cases, we kept the standard when it was
182 observed in the raw data (2,321 instances). In only 43 instances, a standard sequence was either too
183 low in abundance or too frequently mutated to be annotated in either the raw or error-corrected
184 datasets.

185 In order to assess potential biases introduced by library preparation we prepared control libraries
186 containing only the pool of synthetic standards (from the master stock). The libraries were
187 generated in the same manner as the described MAF protocol, with the exception that in the first
188 PCR step, instead of using a multiplex forward primer set, a single forward primer targeting the
189 conserved 5' non-coding region (singleplex-PCR) was used. Ig-seq on these samples allowed us to
190 establish a baseline for pipetting accuracy by comparing the obtained standard frequencies from the
191 singleplex-PCR against the expected frequencies based on our pooling scheme: this yielded an R^2
192 of 0.88 and average mean squared error (MSE) of $0.29 \pm 0.02\%$ (**Figure S1A**). These values
193 indicate that only small systematic deviations occurred, most likely due to minor pipetting error.
194 Next, we compared standard frequencies (expected relative concentration) with frequencies
195 generated in our previous multiplex-PCR libraries, both with and without MAF correction (**Figure**
196 **3C**). On uncorrected data, the multiplex-PCR libraries achieved an R^2 of 0.84 with an average MSE
197 of $0.34 \pm 0.06\%$, which is significantly worse than the value obtained by singleplex-PCR (Student's
198 t-test $p = 0.008$). After error and bias correction on these same datasets, the correlation improved to
199 an R^2 of 0.89 and an average MSE of $0.28 \pm 0.08\%$. While MAF-corrected MSE values show no
200 significant difference to the singleplex-PCR libraries (Student's t-test $p = 0.49$), they do highlight a
201 significant improvement over the uncorrected data (paired Student's t-test, $p = 0.0007$).

202 **Impact of MAF error correction on human B cell repertoires**

203 Next, we analyzed the impact of MAF on the BCR repertoires of B cells isolated from the
204 peripheral blood of three healthy donors. We used flow cytometry sorting and a gating strategy to
205 select for $CD27^- IgM^+$ (naïve) and $CD27^+ IgG^+$ (memory) B cell populations (**Figure 4A**). Across
206 all donors, the fraction of $CD19^+$ peripheral blood B cells was 16-29% for $CD27^- IgM^+$ and 5-9%
207 for $CD27^+ IgG^+$. Importantly, each donor population was split into 4-5 separate aliquots containing
208 200,000 cells each (cellular replicates) prior to cell lysis. Total RNA was extracted and RT for
209 cDNA synthesis was performed independently in order to prevent the mixing of transcripts across
210 cellular replicates. The cDNA of synthetic standards (from the master stock) was then mixed with
211 the B cell cDNA, and corresponding molecular quantities were measured by ddPCR (**Figure 4B**,
212 **Table S3**). All cDNA libraries were then processed into libraries using the MAF protocol (**Figure**
213 **2A**) and subjected to Ig-seq.

214 A simple global analysis of Ig-seq data revealed that diversity measurements were dramatically
215 exaggerated, as the number of unique antibody sequence variants obtained from the raw,

216 uncorrected data often exceeded both the number of cells and the number of total cDNA transcripts
217 in a given aliquot. Following error correction by MAF, the variant count returned to ranges that are
218 physically possible, thereby highlighting the importance of proper error correction and quality
219 control when globally determining repertoire diversity (**Figure 4B**). We further examined the
220 influence of erroneous variants on CDR3 clonotype analysis. In order to identify clonotypes, we
221 used hierarchical clustering (30) based on sequences sharing the following features: identical IGHV
222 and IGHJ gene segment usage, identical CDR3 length, and a CDR3 a.a. similarity of at least 80% to
223 one other sequence in the given clonotype. When performing such an analysis on uncorrected data,
224 clonotypes contained an artificially high number of distinct clones (**Figure 4C**, left tree). Here, the
225 IgG-derived clonotype with the consensus CDR3 of 'CARAAGSQYYMDVW' (from the same
226 sample and IGHV gene segment used by the standard shown in **Figure 3D**) contains 249 unique nt
227 variants and 70 unique a.a. sequences. After MAF-based error correction, only 15 nt variants and 6
228 distinct CDR3 a.a. sequences remained. It is worthy to note that although both the standard
229 sequence (**Figure 3D**) and the biological clonotype (**Figure 4C**) had a large number of CDR3 a.a.
230 variants in uncorrected data, after error correction only the biological clonotype retained multiple
231 a.a. variants, suggesting these may be true variants generated *in vivo* by SHM. This general trend of
232 each IgG memory B cell-derived clonotype to contain more variants relative to antigen
233 inexperienced IgM-expressing B cells was clear across our biological data sets (**Figure 5A**).

234 **Clonal diversity measurements of human B cell repertoires after error and bias correction**

235 After establishing the value of performing MAF error correction on biological repertoires, we next
236 focused on determining the clonotype diversity present in each B cell sample. First, we determined
237 the overlap of clonotypes present in each cellular replicate (**Figure 6A**). Notably, we observed an
238 overlap of several clonotypes in the CD27⁺IgM⁺ subset; this overlap was unexpected given that this
239 subset should be highly enriched for naïve B cells, which by definition are not antigen experienced
240 or clonally expanded, and should therefore be mostly unique (not present in multiple replicates).
241 For each donor in the CD27⁺IgM⁺ subset, 1- 2% of all clonotypes were present in at least one other
242 cellular replicate. In the CD27⁺IgG⁺ subset, clonotypes shared between at least two cellular
243 replicates were nearly tenfold more frequent (12-15%), which was expected given that this
244 population is comprised of antigen experienced, clonally expanded memory B cells. Another
245 observation discordant with the expected naïve B cell properties of the CD27⁺IgM⁺ subset was that
246 overlapping clonotypes (in donors 1 and 3) were significantly more likely to have acquired
247 mutations (**Figure 6B**), which are not a typical feature of naïve B cells. In comparison, over 90% of
248 all CD27⁺IgG⁺ (overlapping and non-overlapping) clonotypes possessed at least one SHM, an
249 expected observation in a memory subset.

250 The high amount of overlap within the CD27⁺IgG⁺ B cell replicates of each donor allowed us to use
251 established population diversity estimation techniques to calculate clonal diversity (31). Rarefaction
252 curves were generated and estimates were extrapolated as a function of real and predicted cellular
253 replicates (**Figure 6C**). The asymptote was determined by the standard form of the Chao2 estimator
254 and yielded the following values for clonotype numbers: donor 1 = 164,268 ± 2,365, donor 2 =
255 38,034 ± 1,302, and donor 3 = 76,904 ± 1,409. Since the 95% confidence intervals for the three
256 donors did not overlap, we could also infer that the size of each donor's repertoire at the collection
257 time point was significantly different. This analysis indicates that we would need to sample at least
258 tenfold more cellular replicates in order to observe > 90% of all clonotypes; however the first five

259 samples analyzed here were sufficient to observe > 25% of the clonotypic memory repertoire. We
260 also generated rarefaction and extrapolation curves rescaled to the RID count (**Figures S2A and**
261 **S2B**). In the case of the CD27⁺IgM⁺ repertoire data, while asymptotic curves could be generated, a
262 diversity estimation is impractical. This is because plotting the observed numbers of newly
263 discovered clonotypes for each additional RID and donor shows that the number of newly
264 discovered clonotypes in the CD27⁺IgM⁺ dataset continues to grow over the observed range,
265 whereas the number of new clonotypes starts to converge at approximately 20,000 RIDs for
266 CD27⁺IgG⁺ repertoires (**Figure S2C**).

267 **Divergent features of CD27⁺IgM⁺ and CD27⁺IgG⁺ repertoires**

268 After pooling all clonotypes (expanded and unique to a single cellular replicate) for each donor, we
269 globally characterized sequences of the naïve CD27⁺IgM⁺ and memory CD27⁺IgG⁺ subsets. First,
270 we determined the SHM count (nt) of each clone with respect to its nearest germline IGHV and
271 IGHJ gene segment sequence. The median values of SHM for the CD27⁺IgM⁺ repertoires were
272 zero, which was to be expected for a naïve B cell subset. In contrast, the median values for
273 CD27⁺IgG⁺ repertoires were 20-24 mutations per clone (**Figure 5B**). It is widely appreciated that
274 human heavy chain CDR3 sequences are much longer than their murine counterparts, which we
275 also observed here, with a slight (but consistent across donors) variation between naïve and memory
276 B cells (**Figure 5C**). Interestingly, the IGHV gene segment family usage correlated with B cell
277 subset. The CD27⁺IgG⁺ repertoires across all donors were relatively enriched for IGHV1 and
278 IGHV3 gene segment family members, whereas the relative share of the IGHV4 gene segment
279 family was larger in the CD27⁺IgM⁺ repertoires (**Figure 7A**). We validated these observations
280 quantitatively using linear discriminant analysis (LDA) fitted on the centered log ratio (CLR)-
281 transformed frequencies of each cellular replicate (**Figure S3**). The LDA classifier was fit on
282 different splits of the data (based on two of the donors) and used to predict a holdout set (based on
283 the remaining donor). This showed that the fitted classifier in each instance is highly predictive of
284 the remaining aliquots and that prediction is robustly driven by the relative abundance of IGHV4
285 segment family usage in the CD27⁺IgM⁺ repertoires and the IGHV1, 2, and 3 families in the
286 CD27⁺IgG⁺ repertoires (**Figure S3A-C**). Next, we utilized LDA to perform dimensionality
287 reduction of all data points to into a single one-dimensional axis; again, the most important
288 components were the relative abundance of IGHV3 and IGHV4 gene segment families (**Figure**
289 **S3D**).

290 Next, we leveraged the ability of our reverse primer to distinguish among IgG subclasses (**Figure**
291 **7B**). The majority of sequences mapped either to IgG1 (40-66%), IgG2 (23%-36%), or IgG3 (23%-
292 36%), whereas IgG4 sequences were extremely rare, observed solely in donor 3 (0.3%). Finally, we
293 compared the CD27⁺IgM⁺ and CD27⁺IgG⁺ repertoires of each donor to determine the clonotype
294 overlap of each B cell subset and isotype. Strikingly, the observed overlap was very small, with
295 only 269 shared clonotypes for donor 1, 30 shared clonotypes for donor 2, and 215 clonotypes for
296 donor 3 (**Figure 7C**). In each donor, these represented less than 0.5% of identified clonotypes.
297 Closer examination revealed that clonotypes shared between the CD27⁺IgM⁺ and CD27⁺IgG⁺
298 subsets were also significantly enriched for intraclonal variants (SHM in CDR3) in one of the two
299 populations (**Figure 7D and Figure S4**). Furthermore, we could see that clonotypes with multiple
300 IgM variants were also shared specifically among IgM cellular replicates (**Figure 7D**, contingency
301 tables). This intraclonal variant bias was not limited to heavy chain isotype: IgG clonotypes with \geq

302 5 intraclonal variants also exhibited subclass composition skewing toward the IgG1/2 or IgG1/3
303 axis, but rarely at proportions similar to the overall IgG subclass distribution (**Figure 7E**).

304

305 **DISCUSSION**

306 Ig-seq is becoming an essential tool for the quantitative analysis of antibody repertoire diversity and
307 distribution. However, similar to other areas of high-throughput sequencing, Ig-seq also suffers
308 from technical errors and bias; thus, standardized experimental and analytical methods that increase
309 the validity of immunological interpretations must be developed. Here, we establish a
310 comprehensive set of synthetic standards which, when combined with UID-labeling and MAF-
311 based error and bias correction, results in highly accurate antibody repertoire data. By applying this
312 approach to human B cell subsets, we gain unique insights into repertoire features such as clonal
313 diversity, germline gene usage, SHM, and clonal history.

314 The synthetic standards developed here allowed quantitative interrogation of several accuracy-
315 related features in Ig-seq. One major observation was that raw uncorrected data has a high number
316 of erroneous variants, found both within the clonotype-defining CDR3 and across the entire VDJ
317 region (**Figure 3A, B**). The number of false-positive variants correlated with the abundance of each
318 standard; this is of particular concern because high frequency, clonally expanded B cells are often
319 correlated with antigen specificity and thus important for biological interpretations (32). However,
320 when applying our MAF-error correction protocol, we were able to remove nearly all erroneous
321 CDR3 and VDJ variants (94-97%); this correction was robust even for high-frequency standards
322 where the number of erroneous variants was especially high. Errors not removed by MAF could
323 potentially be addressed with more stringent filtering criteria (e.g. read number cutoffs); however,
324 this may come at the cost of reducing overall dataset size and removing legitimate intraclonal
325 variants in biological samples.

326 Another aspect we quantified with our standard pool was the impact of multiplex primer sets, which
327 have been shown to introduce substantial bias during library preparation (18, 33). By designing our
328 standards with a 5' conserved singleplex region (**Figure 1A**), we were able to directly compare Ig-
329 seq data from libraries (on the same master stock) prepared by singleplex-PCR vs. multiplex-PCR.
330 Our newly designed FR1-targeting multiplex primer set (**Figure 2B, C**) demonstrated a relatively
331 strong correlation with singleplex-PCR data ($R^2 = 0.84$) (**Figure 3C**). However, by performing
332 multiple MAF error and bias correction steps, the correlation was improved by an additional 7%
333 (**Figure 3D**). The remaining variability does not appear to be restricted to a particular IGHV-gene
334 family, indicating there is little systematic bias with respect to homologous sequences within the
335 standard pool. MAF therefore represents an essential step in eliminating technical artifacts from
336 human Ig-seq workflows, as it is able to generate data that closely mirrors that of the original
337 sample. In future applications, these synthetic standards could be a critical asset for evaluating
338 newly designed primer sets, library preparation protocols, or implementing new error and bias
339 correction pipelines.

340 Having established a comprehensive set of synthetic standards and a validated error and bias
341 correction pipeline, we were able to perform several analyses on human B cell repertoires with
342 greater confidence in the accuracy and quantitative resolution of the Ig-seq data. A simple approach

343 to estimating repertoire diversity is to calculate the number of unique antibody sequences as a
344 fraction of total transcript (cDNA) input. However, performing this analysis on our samples
345 suggests that the CD27⁺IgG⁺ memory B cell compartment is significantly more diverse than the
346 naïve CD27⁺IgM⁺ B cells (82% vs. 35% unique nt variants, respectively, averaged across donors
347 and cellular aliquots, Student's t-test $p < 10^{-4}$). This is potentially due to sample size variability with
348 respect to the number of transcripts and our ability to oversample smaller libraries. Critically, when
349 using bulk-sorted cells with UID labeling, it is not possible to discriminate between transcript
350 copies that are identical because they came from the same lysed cell, and those which are identical
351 because they represent two distinct, clonally related B cells. Thus, by biological subsampling
352 through cellular replicates, we ensured that clonotypes observed in multiple samples must come
353 from distinct, clonally related B cells, thereby providing an effective solution for estimating clonal
354 diversity.

355 Applying computational approaches from ecology (31) to our biological subsampling strategy, we
356 attempted to estimate the number of unique clonotypes in a given antibody repertoire. The CD27⁻
357 IgM⁺ B-cell subset did not show substantial clonotype overlap among cellular aliquots (**Figure 6A**).
358 As it is commonly assumed (and typically the case as shown in **Figure 5A**) that each newly
359 generated B cell is a unique clone, the size of the naïve repertoire in the human peripheral blood
360 would be equal to the total number of naïve B cells, in the range of 10 to 30 million. While it is
361 improbable to sample this subset in its entirety, and its diversity is also too high to estimate based
362 on the cell numbers obtained here, our observations are consistent with this model, since each
363 additional cellular replicate produced overwhelmingly unique sequences. One donor did show an
364 unexpected presence of overlapping sequences (shared clonotypes) across IgM cellular replicates
365 (**Figure 6A**, donor 1); these clonotype sequences were significantly enriched for SHM (**Figure 6B**),
366 suggesting the possible presence of an antigen-experienced B cell subset within CD27⁻IgM⁺
367 population, and highlighting the need for improved characterization of the heterogeneity within
368 circulating human B cell subsets. In the CD27⁺IgG⁺ B cell subset, we observed substantially more
369 overlap across cellular replicates, which was expected given that memory-enriched B cells would
370 have experienced antigen and undergone clonal expansion. By extrapolating the numbers of
371 additional uniquely observed clonotypes with each subsequent cellular aliquot, we were able to
372 predict the clonotype size of the peripheral CD27⁺IgG⁺ B-cell repertoire to be on the order of 10^5
373 (**Figure 6C**). Indeed, rough estimates of the number of antigen-specific clonotypes generated by a
374 single immune response (≈ 100 , a number in line with what has been described regarding serum
375 antibody clonotypes (34)) and the number of structurally distinct pathogens against which an
376 individual has mounted a response (≈ 1000 , a generous estimate given work showing that
377 worldwide, individuals have on average a serological history against less than 100 viral species
378 (35)) suggest that a memory repertoire of this size could reasonably protect against latent infection
379 and/or subsequent antigen encounter.

380 We observed a clear shift in IGHV segment family usage from the naïve to the memory BCR
381 repertoire (**Figure 7A**). Consistently observing this reshaping in three independent healthy donors,
382 and comparing to our standard controls to exclude the possibility of biased amplification, we can
383 conclude that it is a genuine phenomenon. Relatively more abundant IGHV1 and less abundant
384 IGHV4 segment usage in IgG memory B cells has been previously observed in one three-donor
385 cohort (1) but not in another which pooled sequences from both class-switched and IgM-memory

386 cells (36), underscoring the importance of experimental design and accurate bias correction in
387 antibody repertoire analysis.

388 Our Ig-seq workflow also allowed us to unambiguously assign IgG antibody sequences to their
389 appropriate subclass, offering further insight into patterns of class-switch recombination present in
390 memory-enriched B cells. While plasma cell-secreted IgG proteins in human serum are present at
391 ratios of approximately 14:8:1:1 (for IgG1:2:3:4, respectively (37)), CD27⁺IgG⁺ B cells showed a
392 distribution of approximately 5:3:1 for IgG1:2:3, with a nearly complete absence of IgG4 (**Figure**
393 **7B**). Cole et al. similarly observed a lack of IgG4 heavy chains in a single donor but described an
394 enrichment of IgG2 relative to IgG1 and IgG3 (24). The abundance of IgG3⁺ B cells relative to its
395 presence in the serum seen here indicates IgG3 may play a more important role in maintaining the
396 reactive memory response compared to the protective memory response provided by serum IgG.
397 Notably, the IgG3 locus is the most proximal, and thereby the most plastic of the human IgG
398 subclasses; that is, an IgG3⁺ B cell still retains the capacity to class-switch to any of the remaining
399 three IgG subclasses, whereas IgG1, IgG2, and IgG4 cannot return to any of the previous states.
400 Similar to these findings, a flow cytometry-based investigation has also found healthy human
401 donors to have low frequencies of IgG4-expressing circulating memory B cells (38).

402 With new daily production and relatively rapid turnover of naïve B cells, it was not unexpected to
403 see little overlap of clonotypes between the intradonor CD27⁺IgM⁺ and CD27⁺IgG⁺ populations
404 (**Figure 7C**). An interesting finding was that for clonotypes present in both B cell subsets,
405 intraclonal variation was largely restricted to one of the two isotypes. Assuming that clonotype
406 overlap among CD27⁺IgM⁺ cellular replicates represents the presence of antigen-experienced,
407 clonally expanded B cells, this suggests that the antigen specificity of antibody variable domains
408 may to some extent be influenced by the downstream constant regions, which has been observed
409 functionally for small cohorts of human and murine IgG and IgA antibodies (39). Notably, we also
410 observe similar clonal restriction within IgG clonotypes with respect to heavy chain subclass
411 (**Figure 7E**). This may be driven by the type of antigen and the nature of the elicited immune
412 response, or governed by physical constraints as we suggest for the differences between the IgM
413 and IgG repertoires. A larger scale functional study, including IgM sequences, could provide crucial
414 support for this model, which would shed new light on the role of Ig isotypes and subclasses on B
415 cells in the post-antigen encounter setting.

416 **FIGURE LEGENDS**

417
418 **Figure 1.** A comprehensive set of human synthetic spike-in standards for Ig-seq. **(A)** Schematic
419 showing the prototypical spike-in with the following regions (5' to 3'): a conserved non-coding
420 region, ATG start codon, IGHV region with FR1 specific for multiplex-PCR, IGHJ regions, non-
421 coding synthetic sequence identifier (specific for ddPCR probes), and downstream heavy chain
422 constant domain sequences (IGHG, IGHM, IGHA) containing primer binding sites used for cDNA
423 synthesis. Each spike-in contains a complete VDJ open reading frame, including nucleotides
424 upstream of the ATG start codon, and downstream constant domain sequences containing primer
425 binding sites used for sample cDNA synthesis. **(B)** Pairwise comparisons based on a.a. Levenshtein
426 edit distance of all 85 standard CDR3 sequences. The germline IGHV and IGHJ segment family
427 usage and IgG subclasses are denoted. 39 spike-ins contain rationally designed nt SHM (black
428 circles) across the IGHV regions.

429

430 **Figure 2.** Library preparation of immunoglobulin (Ig) heavy chain genes for high-throughput
431 sequencing (Ig-seq) using molecular amplification fingerprinting (MAF). **(A)** In step 1, reverse
432 transcription (RT) is performed to generate first-strand cDNA with a gene-specific (IgM or IgG)
433 primer which includes a unique reverse molecular identifier (RID) and partial Illumina adapter (IA)
434 region. This results in single-molecule labeling of each cDNA with an RID. In step 2, several cycles
435 of multiplex-PCR are performed using a forward primer set with gene-specific regions targeting
436 heavy chain variable (V_H) framework region 1 (FR1), with overhang regions comprised of a
437 forward unique molecular identifier (FID) and partial IA. In step 3, singleplex-PCR is used to
438 extend the partial IAs. The result (Step 4) is the generation of antibody amplicons with FID, RIDs,
439 and full IA ready for Ig-seq and subsequent MAF-based error and bias correction. **(B)** List of
440 oligonucleotides sequences annealing to the V_H FR1 used in multiplex-PCR (Step 2) of the MAF
441 library preparation protocol. The nearest germline IGHV segment(s) likely to be amplified by the
442 respective primer are listed in the rightmost column. **(C)** The estimated melting temperature
443 distribution of the V_H FR1 forward primer set.

444

445 **Figure 3.** Synthetic standards used to validate performance of error and bias correction of Ig-seq
446 data by MAF. **(A)** The number of erroneous CDR3 sequences (at least one a.a. difference from the
447 correct CDR3) per 100,000 reads is plotted against the relative concentration of each standard (from
448 a master stock). Color-coded diamonds correspond to germline IGHV segment family of the
449 respective standard and show the number of erroneous variants in uncorrected (raw) data; gray
450 diamonds indicate the number of variants remaining after MAF error correction. **(B)** The number of
451 erroneous VDJ variants derived from each standard was calculated by finding all variants that
452 carried the correct CDR3 a.a. sequence, but differed by at least one nt across the entire VDJ region.
453 Colored diamonds represent uncorrected data; gray diamonds indicate variants remaining after
454 MAF error correction. **(C)** Sequencing bias introduced by multiplex-PCR using the FR1 primer set
455 was assessed by plotting the measured frequencies of each standard against its relative
456 concentration (from a master stock). Dashed line represents a bias-free ideal scenario ($R^2 = 1$). The
457 left and right plots show observed frequencies before and after MAF bias correction, respectively.
458 **(D)** Phylogenetic trees visualizing the CDR3 a.a. variants present for a selected standard with the
459 CDR3 a.a. sequence *CARGINGERALEW* and IGHV1-8 and IGHJ1 segment usage. Prior to error
460 correction, 39 erroneous CDR3 a.a. variants (branches) and 218 VDJ nt variants (black circles)
461 were observed. Following MAF error correction, only the original correct CDR3 a.a. and two VDJ
462 nt variants remain. The Ig-seq data sets used in A-C consisted of ~300,000 preprocessed full-length
463 antibody reads from each of the synthetic spike-in only samples. IgG1_D1 dataset was used for
464 panel **(D)** (see **Table S3**).

465

466 **Figure 4.** Ig-seq analysis of human naïve ($CD27^-IgM^+$) and memory ($CD27^+IgG^+$) B cells. **(A)** The
467 flow cytometric workflow for isolating $CD27^-IgM^+$ and $CD27^+IgG^+$ B cells from peripheral blood.
468 Boxed-in values indicate the frequency of each sorted subset as a percentage of the total B cell
469 ($CD19^+$) population from each of three donors (1-3 from left to right). **(B)** Experimental and Ig-seq
470 based quantitation of antibody diversity; points represent cDNA molecule counts (using ddPCR) or
471 unique reads (before and after MAF error correction) from cellular replicates (with mean and
472 standard deviation shown) isolated from each donor and B cell subset. Unique read counts were
473 based on the VDJ nt sequence. Dashed line represents the number of B cells isolated per cellular
474 replicate (2×10^5 cells). **(C)** Phylogenetic trees illustrating CDR3 a.a. and nt variants present for the

475 selected clonotype with the consensus CDR3 sequence *CARAAGSQYYYMDVW* and IGHV1-8 to
476 IGHJ1 recombination. Prior to error correction, 70 erroneous CDR3 a.a. variants and 249 VDJ nt
477 variants (black circles) were observed. Following MAF error correction, only 6 CDR3 a.a. and 15
478 VDJ nt variants remain. The Ig-seq data sets used in **(B)** are described in **Table S3**; IgG1_D1 was
479 used for the tree in panel **(C)**.

480

481 **Figure 5.** Ig-seq analysis of molecular features highlight global differences between naïve (CD27⁻
482 IgM⁺) and memory (CD27⁺IgG⁺) B cells. **(A)** Clonotype size (calculated as the total number of
483 variants within a clonotype) for naïve IgM (blue lines) and memory IgG (red lines) B cells. Each
484 pair of lines represents a single donor. **(B)** Graph showing the distribution of average SHM
485 frequencies for VDJ nt variants per clonotype. The CD27⁻IgM⁺ B cell repertoires (blue lines) have a
486 median SHM value of zero, whereas only a small fraction of clonotypes (approximately 8%)
487 contain an average of one or more mutations. CD27⁺IgG⁺ repertoires (red lines) have a median of
488 20 to 24 SHM per nt variant within each clonotype. **(C)** CDR3 a.a. length distribution across
489 clonotypes from naïve (blue lines) and memory (red lines) B cell subsets.

490

491 **Figure 6.** Clonotype diversity analysis across cellular replicates of naïve (CD27⁻IgM⁺) and memory
492 (CD27⁺IgG⁺) B cells. **(A)** Venn diagrams show the presence of clonotypes (80% CDR3 a.a.
493 similarity to least one clone in the cluster, same CDR3 a.a. length, same IGHV and IGHJ gene
494 segment usage) across cellular replicates (2×10^5 cells each) from each donor B cell subset. **(B)** Bar
495 graph showing the fraction of clonotypes containing at least one variant with at least one nt SHM.
496 Blue and red bars indicate clonotypes identified either in only one aliquot (unique) or in several
497 aliquots (shared), respectively. The p-values represent significance using Fisher's exact test. **(C)**
498 Species accumulation curves for CD27⁺IgG⁺ B cells: the number of newly discovered clonotypes
499 from each additional cellular replicate (black circles) is plotted. Extrapolating the observed overlap
500 provides an estimate for the total number of distinct clonotypes (Chao2 estimator: $D1 = 164,268 \pm$
501 $2,365$; $D2 = 38,034 \pm 1,302$; $D3 = 76,904 \pm 1,409$) and the approximate amount of cellular
502 replicates needed to discover all clonotypes present in the peripheral blood CD27⁺ IgG⁺ population.

503

504 **Figure 7.** Genetic features of naïve and IgG memory BCR repertoires. **(A)** Block map shows IGHV
505 gene segment usage sorted by family (color-coded) across donors and B cell subset; blocks are
506 normalized to the total number of clonotypes within each group. **(B)** IgG subclass usage in
507 CD27⁺IgG⁺ donor repertoires. IgG₄ sequences (dark blue bars) were virtually absent among three
508 donors; a small fraction of sequences could not be unambiguously mapped based on the sequencing
509 read (gray bars). **(C)** Venn diagrams showing the overlap of clonotypes (80% CDR-H3 amino acid
510 similarity to least one clone in the cluster, same CDR3 a.a. length, same IGHV and IGHJ gene
511 segment usage) between the naïve (CD27⁻IgM⁺) and memory (CD27⁺IgG⁺) BCR heavy chain
512 repertoire in each of three donors. **(D)** Each plot shows the clonal composition of each shared
513 clonotype (from panel **(C)**) in terms of its IgG and IgM intraclonal variants. Red triangles indicate
514 clonotypes found in multiple IgM cellular aliquots; blue triangles show clonotypes which could
515 only be found in one IgM cellular aliquot. The total number of clonotypes found are depicted in the
516 corresponding contingency table. Fisher's exact test was used to quantitatively analyze enrichment
517 of expanded IgM clonotypes in the shared IgM/IgG subset. **(E)** Ternary plots comprised of three
518 axes representing the IgG1, IgG2, and IgG3 isotype subclasses. The relative subclass composition
519 of intraclonal variants per IgG clonotype (each represented by a circle colored according to the
520 number of variants belonging to that clonotype) is depicted by the position of the circle within the

521 triangular space. Red circles represent the average subclass composition for all IgG variants of each
522 donor (cf. **Figure 6C**).

523

524

525

526

527 **METHODS**

528

529 **Preparation of spike-in master stocks**

530

531 The spike-in standards were ordered from GeneArt (Invitrogen) in the form of plasmids. Each
532 spike-in sequence contained a T7 promoter for *in vitro* transcription. Approximately 1.5 µg of each
533 plasmid was digested with 10 U of EcoRV-HF (New England BioLabs) and purified with DNA-
534 binding magnetic beads (SPRI select, Beckman Coulter). Approximately 1 µg of the digested
535 plasmid was then used for *in vitro* transcription (MEGAscript T7 Transcription Kit, ThermoFisher
536 Scientific). RNA was purified by lithium-chloride precipitation, eluted (TE with 1U/µl RiboLock)
537 and aliquoted. The final concentration was then determined with the TapeStation (Agilent
538 Technologies).

539 The spike-in RNA obtained this way was reverse transcribed with the Maxima Reverse
540 Transcriptase kit (Thermofisher Scientific). 500 ng mRNA was mixed with 20 pmol of IgM reverse
541 primer and 3 µl dNTP-mix (10 mM each) and was then filled up to 14.5 µl with water. The
542 reaction-mix was incubated at 65 °C for 5 minutes. 4 µl of 5x RT-buffer, 0.5 µl (20 U) RiboLock
543 and 1 µl Maxima reverse transcriptase (200 U) were then added. The resulting reaction mix was
544 incubated for 35 minutes at 55 °C, followed by a termination step at 85 °C for 5 minutes. 2.5 µl of
545 RNase A (Thermofisher Scientific) was added and the mix was again incubated at 60 °C for 30
546 minutes. The resulting cDNA was purified with SPRI Select magnetic beads and eluted in nuclease-
547 free water. The concentration of each cDNA reaction was determined afterwards with the Fragment
548 Analyzer and pooled according to Table S1. The exact concentration of the pooled spike-ins was
549 determined by ddPCR with dilutions of the pool ranging from 10⁻³ to 10⁻⁶. The measured spike-in
550 pool was afterwards diluted to a final storage concentration of 250,000 transcripts per µl.

551

552 **Transcript quantitation by ddPCR**

553

554 Quantifying cDNA and PCR products by ddPCR was conducted for all measurements in the
555 following way: A dilution series with 3 or 4 points was prepared. Droplets were generated with
556 BioRad's droplet generator using 12.25 µl of ddPCR Supermix (BioRad) combined with 10 µl of
557 the diluted sample, 25 pmol of the biological ddPCR probe (SF_21), 25 pmol of the spike-in
558 specific ddPCR probe (TAK_499), 22.5 pmol of the forward (SF_63) and reverse ddPCR primer
559 (TAK_522) and 55 µl of droplet generation oil (BioRad). Droplets were then transferred to a 96-
560 well reaction plate, which was heat sealed with easy pierce foil (VWR International). Then a PCR
561 reaction was performed using the following conditions: 95°C for 10 min; 45 cycles of 94 °C for
562 30s, 53 °C for 30s, 64 °C for 1 min; 98°C for 10 min; and holding at 4 °C. After the PCR step,
563 every 96-well plate was read using BioRad's droplet reader.

564

565 **B-cell isolation, sorting, and lysis**

566

567 Peripheral blood leukocyte-enriched fractions ('buffy coats') were received from the Bern
568 (Switzerland) blood donation center after obtaining the proper informed consent from healthy
569 human donors. Blood samples were diluted 1:3 with sterile PBS and overlaid on Ficoll-Paque
570 PLUS (GE Healthcare) using LeukoSep conical centrifuge tubes (Greiner Bio-One). Peripheral
571 blood mononuclear cells (PBMCs) were harvested after separation for 30 minutes at 400 x g
572 without braking. Successive centrifugation steps were performed to wash the mononuclear cell
573 fraction and remove residual neutrophils and granulocytes. Total B cells were isolated from PBMCs
574 by negative selection with the EasySep Human B Cell Enrichment Kit (STEMCELL Technologies)
575 according to the manufacturer's instructions. The following fluorescently labeled antibodies were
576 used to stain the enriched B-cell fraction prior to sorting by flow cytometry: anti-CD3-APC/Cy7
577 (clone HIT3a BioLegend # 300318), anti-CD14-APC/Cy7 (clone HCD14 BioLegend #325620),
578 anti-CD16-APC/Cy7 (clone 3G8 BioLegend #302018), anti-CD19-BV785 (clone HIB19,
579 BioLegend #302240), anti-CD20-BV650 (clone 2H7, BioLegend #302336), anti-CD27-V450
580 (clone M-T271, BD Horizon #560448), anti-CD24-BV510 (clone ML5, BioLegend #311126), anti-
581 CD38-PC5 (clone LS198-4-3, Beckman Coulter #A07780), anti-IgD-PEcy7 (clone IA6-2,
582 BioLegend #348210), anti-IgG-Alexa Fluor 647 (Jackson ImmunoResearch #109-606-003) anti-
583 IgA-Alexa Fluor 488 (Jackson ImmunoResearch #109-549-011), anti-IgM-PE (clone SA-DA4,
584 eBioscience #12-9998). Cell sorting was performed on a BD FACS Aria III following the gating
585 strategy depicted in Figure 4A. After sorting, isolated fractions were centrifuged 5 minutes at 300 x
586 g, the supernatants were aspirated, and the cell pellets were re-suspended in 1 ml PBS. Recovered
587 cells were hand-counted using a Neubauer hemocytometer, and aliquots containing equal numbers
588 of cells were prepared from the cellular suspension. These aliquots were centrifuged, and the
589 supernatant aspirated. Cell pellets were lysed directly in 200 µl TRI Reagent (Sigma), allowed to
590 dissociate for 5 minutes at room temperature, then frozen on dry ice prior to storage at -80°C.

591

592 **RNA isolation from sorted B-cell populations**

593

594 Immediately prior to use, Phase Lock Gel (PLG) tubes were pelleted at 12000 - 16000 x g in a
595 microcentrifuge for 20 to 30 seconds. Each TRIzol aliquot was then thawed on ice. After thawing
596 and an incubation time of 5 minutes at room temperature, 1 mL of the TRIzol homogenate was
597 transferred to the phase lock tube. 0.2 mL chloroform was added and the tube was shaken
598 vigorously by hand (~15 seconds). After an incubation time of 3 minutes, the phase lock tube was
599 centrifuged at 12'000 x g during 15 minutes at 4 °C. The resulting upper aqueous phase was
600 transferred to a fresh Eppendorf tube, an equal volume of 70% ethanol was added and the solution
601 was purified on a PureLink RNA column according to the manufacturer's instructions (Life
602 Technologies). Finally, RNA was eluted in 25 µl nuclease-free water.

603

604 **Library preparation and NGS on Illumina's MiSeq**

605

606 First-strand cDNA synthesis was carried out using Maxima reverse transcriptase (Life
607 Technologies). The protocol for one reaction is as follows: A 29 µl reaction mix was prepared using
608 up to 2 µg of RNA together with 40 pmol of the respective gene-specific reverse primer (for IgG
609 sequences, 5'-RID-GTTCTGGGAAGTAGTCCTTGACCAG-3' (IgH_10r); for IgM, 5'-RID-
610 ACGAGGGGGAAAAGGGTTGG-3' (CH1_1r)), 2 µl dNTP (10 mM each) and the required
611 amount of nuclease-free water. This mix is incubated for 5 minutes at 65 °C. A master mix of 8 µl
612 5x RT buffer, 1 µl of (20 U) RiboLock and 2 µl of Maxima reverse transcriptase is then prepared

613 and added to the reaction mix. Finally, the mix is incubated for 30 minutes at 50°C and the reaction
614 is terminated by incubating at 85°C for 5 minutes.

615 The obtained cDNA is then cleaned with a left-sided SPRI-Select bead clean up (0.8x) according to
616 the manufacturer's instructions (Beckman coulter) and subsequently measured by ddPCR.

617 Up to 135,000 cDNA transcripts were then pooled together with 12,500 spike-in transcripts.
618 Multiplex-PCR was performed using an equimolar pool of the forward primer mix (Figure 2B) and
619 the reverse primer (TAK_423) targeting the overhang introduced during cDNA synthesis. Due to
620 low cDNA yields, the first PCR was carried out for 20 cycles and the following cycling protocol: 2
621 min at 95 °C; 20 cycles of 98°C for 20s; 60°C for 50; 72°C for 1min; 72°C and then holding at 4°C.
622 PCR reactions were prepared using 15 µl of Kapa HiFi HotStart ReadyMix (KAPA Biosystems), 50
623 pmol of the forward mix and the reverse primer and 9 µl of the cDNA mix. After the first PCR, we
624 again performed a left-sided bead clean-up (0.8x) and measured the PCR product concentration
625 using ddPCR. We use 800,000 transcripts from PCR 1 as input into the adapter extension PCR. For
626 this PCR 25 µl of Kapa HiFi Hotstart ReadyMix was combined with 25 pmol of the forward primer
627 (TAK_424) and 25 pmol of the index primer (TAK_531) as well as the diluted PCR product.
628 Finally, the reaction volume was adjusted to 50 µl by the addition of nuclease-free water.
629 Thermocycling was performed as follows: 95°C for 5 min; 23 cycles of 98°C for 20 s, 65°C for 15
630 s, 72°C for 15 s; 72°C for 5 min; and 4°C indefinitely. Following second-step adapter extension
631 PCR, reactions were cleaned using a double-sided SPRIselect bead cleanup process (0.5x to 0.8x),
632 with an additional ethanol wash and elution in TE buffer.

633 Libraries were then quantified by capillary electrophoresis (Fragment analyzer, Agilent). After
634 quantitation, libraries were pooled accordingly and sequenced on a MiSeq System (Illumina) with
635 the paired-end 2x300bp kit.

636

637 **Bioinformatic pipeline**

638

639 Paired-end fastq files were merged using PandaSeq (40). Afterwards, sequences were filtered for
640 quality and length using the FASTX toolkit (http://hannonlab.cshl.edu/fastx_toolkit/). After the
641 quality trim, sequences were processed with a custom Python script that performed error correction
642 by consensus building on our sequences and RIDs. In order to utilize as many sequencing reads as
643 possible, we required UIDs to have at least 3 reads, but did not remove sequences that only had one
644 UID group mapping to them. VDJ annotation and frequency calculation was then performed by our
645 in-house aligner (18) which was updated with the human reference database downloaded from
646 IMGT. The complete error-correction and alignment pipeline is available under
647 <https://gitlab.ethz.ch/reddy/MAF>.

648

649 **Statistical analysis**

650

651 All statistical and computational analyses following the alignment step were performed in R.
652 Details about specific tests that were used can be found in the results section and in the figure
653 legends. Scripts are available upon request.

654

655 **Data availability**

656

657 In adherence to the data sharing recommendations of the AIRR community our data is publically
658 available in the following repositories: BioProject, BioSample, SRA and GenBank and can be
659 accessed with the accession number PRJNA430091 (BioProject). The exact data processing steps,

660 including software tools and version numbers can be found on zonodo.org under the following doi:
661 10.5281/zenodo.1201416.

662
663 Likewise, the designed spike-in sequences are also stored on GenBank (Accession number
664 MG785894-MG785978).

665

666 **Author contributions**

667 J.M.L., S.F., E.T., and S.T.R. designed experiments. V.C. performed B-cell enrichment, sorting,
668 and mRNA extraction. M.I. and S.F. prepared IgH libraries. J.M.L. designed primer sequences.
669 J.M.L. and A.Z. designed antibody spike-ins. A.Z., S.M., and M.I. conducted preliminary
670 experiments. S.F. was responsible for the bioinformatics pipeline. J.M.L. and S.F. analyzed data
671 and prepared figures. J.M.L., S.F., E.T., and S.T.R. wrote the paper. All authors provided scientific
672 guidance.

673

674 **Acknowledgments**

675

676 We would like to acknowledge the Genomics Facility Basel of ETH Zurich for Illumina sequencing
677 support, in particular E. Burcklen, K. Eschbach and C. Beisel. We also want to thank H.
678 Ruscheweyh for bioinformatic code support.

679

680 **REFERENCES**

681

- 682 1. Wu Y-C, Kipling D, Leong HS, Martin V, Ademokun AA, Dunn-Walters DK. High-throughput immunoglobulin
683 repertoire analysis distinguishes between human IgM memory and switched memory B-cell populations. *Blood* (2010)
684 116(7):1070-8. doi: 10.1182/blood-2010-03-275859. PubMed PMID: 20457872; PubMed Central PMCID:
685 PMCPMC2938129.
- 686 2. Arnaout R, Lee W, Cahill P, Honan T, Sparrow T, Weiland M, et al. High-resolution description of antibody
687 heavy-chain repertoires in humans. *PLoS ONE* (2011) 6(8):e22365. doi: 10.1371/journal.pone.0022365. PubMed PMID:
688 21829618; PubMed Central PMCID: PMCPMC3150326.
- 689 3. DeKosky BJ, Kojima T, Rodin A, Charab W, Ippolito GC, Ellington AD, et al. In-depth determination and
690 analysis of the human paired heavy- and light-chain antibody repertoire. *Nature medicine* (2015) 21(1):86-91. doi:
691 10.1038/nm.3743. PubMed PMID: 25501908.
- 692 4. Robinson WH. Sequencing the functional antibody repertoire--diagnostic and therapeutic discovery. *Nat Rev*
693 *Rheumatol* (2015) 11(3):171-82. doi: 10.1038/nrrheum.2014.220. PubMed PMID: 25536486; PubMed Central PMCID:
694 PMCPMC4382308.
- 695 5. Williams LD, Ofek G, Schätzle S, McDaniel JR, Lu X, Nicely NI, et al. Potent and broad HIV-neutralizing
696 antibodies in memory B cells and plasma. *Science immunology* (2017) 2(7):eaal2200. doi: 10.1126/sciimmunol.aal2200.
697 PubMed PMID: 28783671.
- 698 6. Doria-Rose NA, Schramm CA, Gorman J, Moore PL, Bhiman JN, DeKosky BJ, et al. Developmental pathway
699 for potent V1V2-directed HIV-neutralizing antibodies. *Nature* (2014) 509(7498):55-62. doi: 10.1038/nature13036.
700 PubMed PMID: 24590074; PubMed Central PMCID: PMCPMC4395007.
- 701 7. Zhu J, Ofek G, Yang Y, Zhang B, Louder MK, Lu G, et al. Mining the antibodyome for HIV-1-neutralizing
702 antibodies with next-generation sequencing and phylogenetic pairing of heavy/light chains. *Proceedings of the National*
703 *Academy of Sciences* (2013) 110(16):6470-5. doi: 10.1073/pnas.1219320110. PubMed PMID: 23536288; PubMed
704 Central PMCID: PMCPMC3631616.
- 705 8. Lavinder JJ, Wine Y, Giesecke C, Ippolito GC, Horton AP, Lungu OI, et al. Identification and characterization of
706 the constituent human serum antibodies elicited by vaccination. *Proceedings of the National Academy of Sciences*
707 (2014) 111(6):2259-64. doi: 10.1073/pnas.1317793111. PubMed PMID: 24469811; PubMed Central PMCID:
708 PMCPMC3926051.
- 709 9. Jiang N, He J, Weinstein JA, Penland L, Sasaki S, He X-S, et al. Lineage structure of the human antibody
710 repertoire in response to influenza vaccination. *Science translational medicine* (2013) 5(171):171ra19-ra19. doi:
711 10.1126/scitranslmed.3004794. PubMed PMID: 23390249; PubMed Central PMCID: PMCPMC3699344.

- 712 10. Roskin KM, Simchoni N, Liu Y, Lee J-Y, Seo K, Hoh RA, et al. IgH sequences in common variable immune
713 deficiency reveal altered B cell development and selection. *Science translational medicine* (2015) 7(302):302ra135-
714 302ra135. doi: 10.1126/scitranslmed.aab1216. PubMed PMID: 26311730; PubMed Central PMCID: PMC4584259.
- 715 11. Palanichamy A, Apeltsin L, Kuo TC, Sirota M, Wang S, Pitts SJ, et al. Immunoglobulin class-switched B cells
716 form an active immune axis between CNS and periphery in multiple sclerosis. *Science translational medicine* (2014)
717 6(248):248ra106-248ra106. doi: 10.1126/scitranslmed.3008930. PubMed PMID: 25100740; PubMed Central PMCID:
718 PMC4176763.
- 719 12. Meng W, Zhang B, Schwartz GW, Rosenfeld AM, Ren D, Thome JJC, et al. An atlas of B-cell clonal distribution
720 in the human body. *Nature Biotechnology* (2017) 35(9):879-84. doi: 10.1038/nbt.3942. PubMed PMID: 28829438;
721 PubMed Central PMCID: PMC5679700.
- 722 13. Lee YN, Frugoni F, Dobbs K, Tirosh I, Du L, Ververs FA, et al. Characterization of T and B cell repertoire
723 diversity in patients with RAG deficiency. *Sci Immunol* (2016) 1(6). doi: 10.1126/sciimmunol.aah6109. PubMed PMID:
724 28783691; PubMed Central PMCID: PMC5586490.
- 725 14. Friedensohn S, Khan TA, Reddy ST. Advanced Methodologies in High-Throughput Sequencing of Immune
726 Repertoires. *Trends in biotechnology* (2017) 35(3):203-14. doi: 10.1016/j.tibtech.2016.09.010. PubMed PMID: 28341036.
- 727 15. Georgiou G, Ippolito GC, Beausang J, Busse CE, Wardemann H, Quake SR. The promise and challenge of
728 high-throughput sequencing of the antibody repertoire. *Nature Biotechnology* (2014) 32(2):158-68. doi:
729 10.1038/nbt.2782. PubMed PMID: 24441474; PubMed Central PMCID: PMC4113560.
- 730 16. Jiang L, Schlesinger F, Davis CA, Zhang Y, Li R, Salit M, et al. Synthetic spike-in standards for RNA-seq
731 experiments. *Genome Research* (2011) 21(9):1543-51. doi: 10.1101/gr.121095.111. PubMed PMID: 21816910; PubMed
732 Central PMCID: PMC3166838.
- 733 17. Shugay M, Britanova OV, Merzlyak EM, Turchaninova MA, Mamedov IZ, Tuganbaev TR, et al. Towards error-
734 free profiling of immune repertoires. *Nature Methods* (2014) 11(6):653-5. doi: 10.1038/nmeth.2960. PubMed PMID:
735 24793455.
- 736 18. Khan TA, Friedensohn S, Gorter de Vries AR, Straszewski J, Ruscheweyh H-J, Reddy ST. Accurate and
737 predictive antibody repertoire profiling by molecular amplification fingerprinting. *Science advances* (2016) 2(3):e1501371.
738 doi: 10.1126/sciadv.1501371. PubMed PMID: 26998518; PubMed Central PMCID: PMC4795664.
- 739 19. Kinde I, Wu J, Papadopoulos N, Kinzler KW, Vogelstein B. Detection and quantification of rare mutations with
740 massively parallel sequencing. *Proceedings of the National Academy of Sciences* (2011) 108(23):9530-5. doi:
741 10.1073/pnas.1105422108. PubMed PMID: 21586637; PubMed Central PMCID: PMC3111315.
- 742 20. Shiroguchi K, Jia TZ, Sims PA, Xie XS. Digital RNA sequencing minimizes sequence-dependent bias and
743 amplification noise with optimized single-molecule barcodes. *Proceedings of the National Academy of Sciences* (2012)
744 109(4):1347-52. doi: 10.1073/pnas.1118018109. PubMed PMID: 22232676; PubMed Central PMCID:
745 PMC3268301.
- 746 21. Kivioja T, Vähärautio A, Karlsson K, Bonke M, Enge M, Linnarsson S, et al. Counting absolute numbers of
747 molecules using unique molecular identifiers. *Nature Methods* (2011) 9(1):72-4. doi: 10.1038/nmeth.1778. PubMed
748 PMID: 22101854.
- 749 22. Vollmers C, Sit RV, Weinstein JA, Dekker CL, Quake SR. Genetic measurement of memory B-cell recall using
750 antibody repertoire sequencing. *Proceedings of the National Academy of Sciences* (2013) 110(33):13463-8. doi:
751 10.1073/pnas.1312146110. PubMed PMID: 23898164; PubMed Central PMCID: PMC3746854.
- 752 23. Turchaninova MA, Davydov A, Britanova OV, Shugay M, Bikos V, Egorov ES, et al. High-quality full-length
753 immunoglobulin profiling with unique molecular barcoding. *Nature Protocols* (2016) 11(9):1599-616. doi:
754 10.1038/nprot.2016.093. PubMed PMID: 27490633.
- 755 24. Cole C, Volden R, Dharmadhikari S, Scelfo-Dalbey C, Vollmers C. Highly Accurate Sequencing of Full-Length
756 Immune Repertoire Amplicons Using Tn5-Enabled and Molecular Identifier-Guided Amplicon Assembly. *The Journal of*
757 *Immunology* (2016) 196(6):2902-7. doi: 10.4049/jimmunol.1502563. PubMed PMID: 26856699.
- 758 25. Lefranc M-P, Giudicelli V, Duroux P, Jabado-Michaloud J, Folch G, Aouinti S, et al. IMGT®, the international
759 ImMunoGeneTics information system® 25 years on. *Nucleic acids research* (2015) 43(Database issue):D413-22. doi:
760 10.1093/nar/gku1056. PubMed PMID: 25378316; PubMed Central PMCID: PMC4383898.
- 761 26. Wang C, Liu Y, Xu LT, Jackson KJL, Roskin KM, Pham TD, et al. Effects of aging, cytomegalovirus infection,
762 and EBV infection on human B cell repertoires. *The Journal of Immunology* (2014) 192(2):603-11. doi:
763 10.4049/jimmunol.1301384. PubMed PMID: 24337376; PubMed Central PMCID: PMC3947124.
- 764 27. Menzel U, Greiff V, Khan TA, Haessler U, Hellmann I, Friedensohn S, et al. Comprehensive evaluation and
765 optimization of amplicon library preparation methods for high-throughput antibody sequencing. *PLoS ONE* (2014)
766 9(5):e96727. doi: 10.1371/journal.pone.0096727. PubMed PMID: 24809667; PubMed Central PMCID:
767 PMC4014543.
- 768 28. Greiff V, Menzel U, Haessler U, Cook SC, Friedensohn S, Khan TA, et al. Quantitative assessment of the
769 robustness of next-generation sequencing of antibody variable gene repertoires from immunized mice. *BMC Immunology*
770 (2014) 15(1):40. doi: 10.1186/s12865-014-0040-5. PubMed PMID: 25318652; PubMed Central PMCID:
771 PMC4233042.
- 772 29. Laserson U, Vigneault F, Gadala-Maria D, Yaari G, Uduman M, Vander Heiden JA, et al. High-resolution
773 antibody dynamics of vaccine-induced immune responses. *Proceedings of the National Academy of Sciences* (2014)
774 111(13):4928-33. doi: 10.1073/pnas.1323862111. PubMed PMID: 24639495; PubMed Central PMCID:
775 PMC3977259.
- 776 30. Greiff V, Miho E, Menzel U, Reddy ST. Bioinformatic and Statistical Analysis of Adaptive Immune Repertoires.
777 *Trends in immunology* (2015) 36(11):738-49. doi: 10.1016/j.it.2015.09.006. PubMed PMID: 26508293.

- 778 31. Colwell RK, Chao A, Gotelli NJ, Lin SY, Mao CX, Chazdon RL, et al. Models and estimators linking individual-
779 based and sample-based rarefaction, extrapolation and comparison of assemblages. *Journal of Plant Ecology* (2012)
780 5(1):3-21. doi: 10.1093/jpe/rtr044.
- 781 32. Reddy ST, Ge X, Miklos AE, Hughes RA, Kang SH, Hoi KH, et al. Monoclonal antibodies isolated without
782 screening by analyzing the variable-gene repertoire of plasma cells. *Nature Biotechnology* (2010) 28(9):965-9. doi:
783 10.1038/nbt.1673. PubMed PMID: 20802495.
- 784 33. Carlson CS, Emerson RO, Sherwood AM, Desmarais C, Chung M-W, Parsons JM, et al. Using synthetic
785 templates to design an unbiased multiplex PCR assay. *Nature communications* (2013) 4:2680. doi:
786 10.1038/ncomms3680. PubMed PMID: 24157944.
- 787 34. Lee J, Boutz DR, Chromikova V, Joyce MG, Vollmers C, Leung K, et al. Molecular-level analysis of the serum
788 antibody repertoire in young adults before and after seasonal influenza vaccination. *Nature medicine* (2016)
789 22(12):1456-64. doi: 10.1038/nm.4224. PubMed PMID: 27820605; PubMed Central PMCID: PMC5301914.
- 790 35. Xu GJ, Kula T, Xu Q, Li MZ, Vernon SD, Ndung'u T, et al. Viral immunology. Comprehensive serological
791 profiling of human populations using a synthetic human virome. *Science (New York, NY)* (2015) 348(6239):aaa0698. doi:
792 10.1126/science.aaa0698. PubMed PMID: 26045439; PubMed Central PMCID: PMC4844011.
- 793 36. DeWitt WS, Lindau P, Snyder TM, Sherwood AM, Vignali M, Carlson CS, et al. A Public Database of Memory
794 and Naive B-Cell Receptor Sequences. *PLoS ONE* (2016) 11(8):e0160853. doi: 10.1371/journal.pone.0160853. PubMed
795 PMID: 27513338; PubMed Central PMCID: PMC4981401.
- 796 37. Vidarsson G, Dekkers G, Rispens T. IgG subclasses and allotypes: from structure to effector functions.
797 *Frontiers in immunology* (2014) 5(16):520. doi: 10.3389/fimmu.2014.00520. PubMed PMID: 25368619; PubMed Central
798 PMCID: PMC4202688.
- 799 38. Heeringa JJ, Karim AF, van Laar JAM, Verdijk RM, Paridaens D, van Hagen PM, et al. Expansion of blood
800 IgG4(+) B, TH2, and regulatory T cells in patients with IgG4-related disease. *J Allergy Clin Immunol* (2017). doi:
801 10.1016/j.jaci.2017.07.024. PubMed PMID: 28830675.
- 802 39. Torres M, Casadevall A. The immunoglobulin constant region contributes to affinity and specificity. *Trends*
803 *Immunol* (2008) 29(2):91-7. doi: 10.1016/j.it.2007.11.004. PubMed PMID: 18191616.
- 804 40. Masella AP, Bartram AK, Truszkowski JM, Brown DG, Neufeld JD. PANDAseq: paired-end assembler for
805 illumina sequences. *BMC bioinformatics* (2012) 13(1):31. doi: 10.1186/1471-2105-13-31. PubMed PMID: 22333067;
806 PubMed Central PMCID: PMC3471323.
- 807

Synthetic standards combined with error and bias correction improves the accuracy and quantitative resolution of antibody repertoire sequencing in human naïve and memory B cells

Simon Friedensohn*¹, John M. Lindner*², Vanessa Cornacchione², Mariavittoria Iazeolla², Enkelejda Miho¹, Andreas Zingg¹, Simon Meng¹, Elisabetta Traggiai², and Sai T. Reddy¹

¹*Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland*

²*Novartis Institutes for BioMedical Research, Basel, Switzerland*

*equal contribution

Main Figures

Supplementary Material

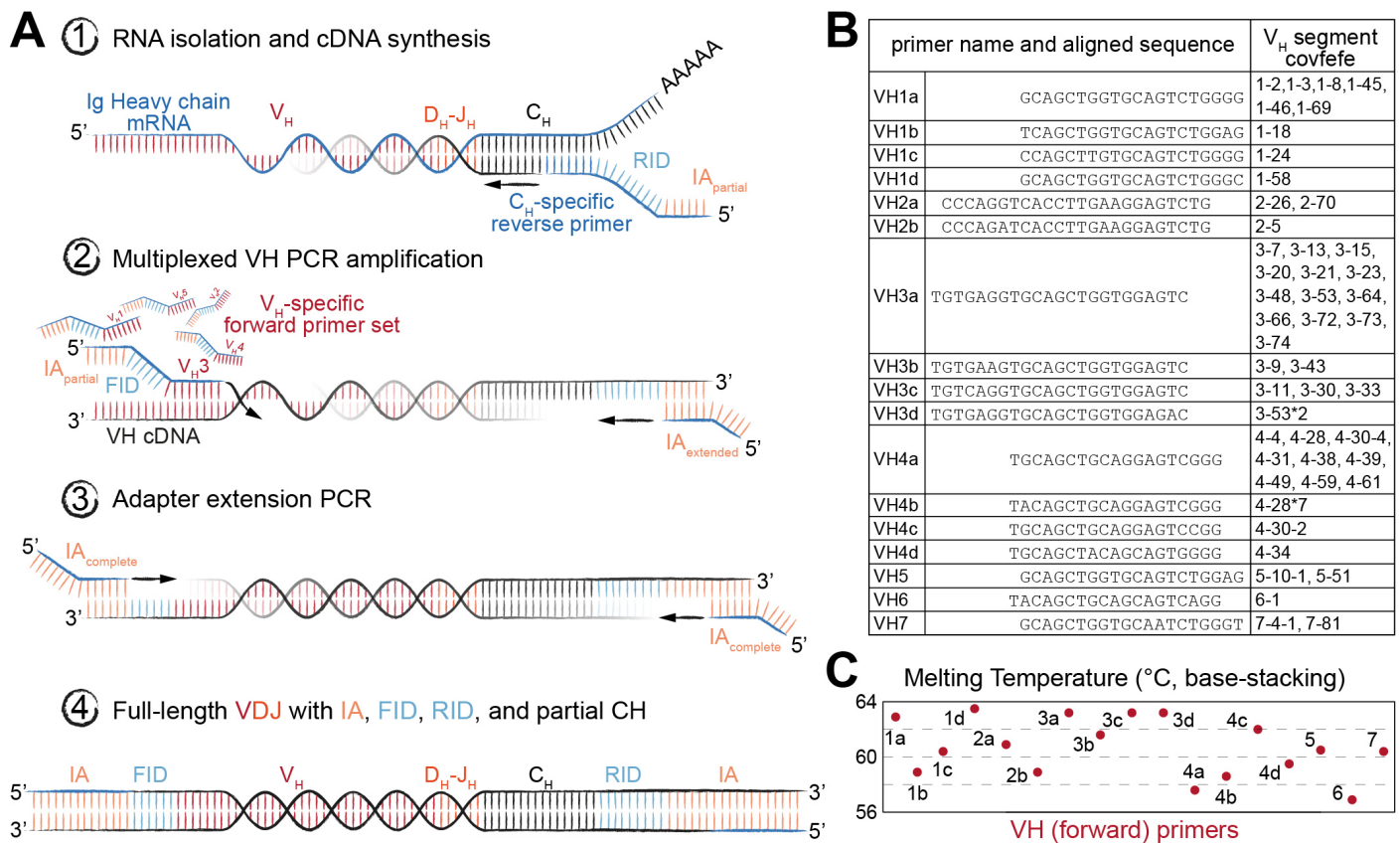


Figure 1. A comprehensive set of human synthetic spike-in standards for Ig-seq. **(A)** Schematic showing the prototypical spike-in with the following regions (5' to 3'): a conserved non-coding region, ATG start codon, IGHV region with FR1 specific for multiplex-PCR, IGHJ regions, non-coding synthetic sequence identifier (specific for ddPCR probes), and downstream heavy chain constant domain sequences (IGHG, IGHM, IGHA) containing primer binding sites used for cDNA synthesis. Each spike-in contains a complete VDJ open reading frame, including nucleotides upstream of the ATG start codon, and downstream constant domain sequences containing primer binding sites used for sample cDNA synthesis. **(B)** Pairwise comparisons based on a.a. Levenshtein edit distance of all 85 standard CDR3 sequences. The germline IGHV and IGHJ segment family usage and IgG subclasses are denoted. 39 spike-ins contain rationally designed nt SHM (black circles) across the IGHV regions.

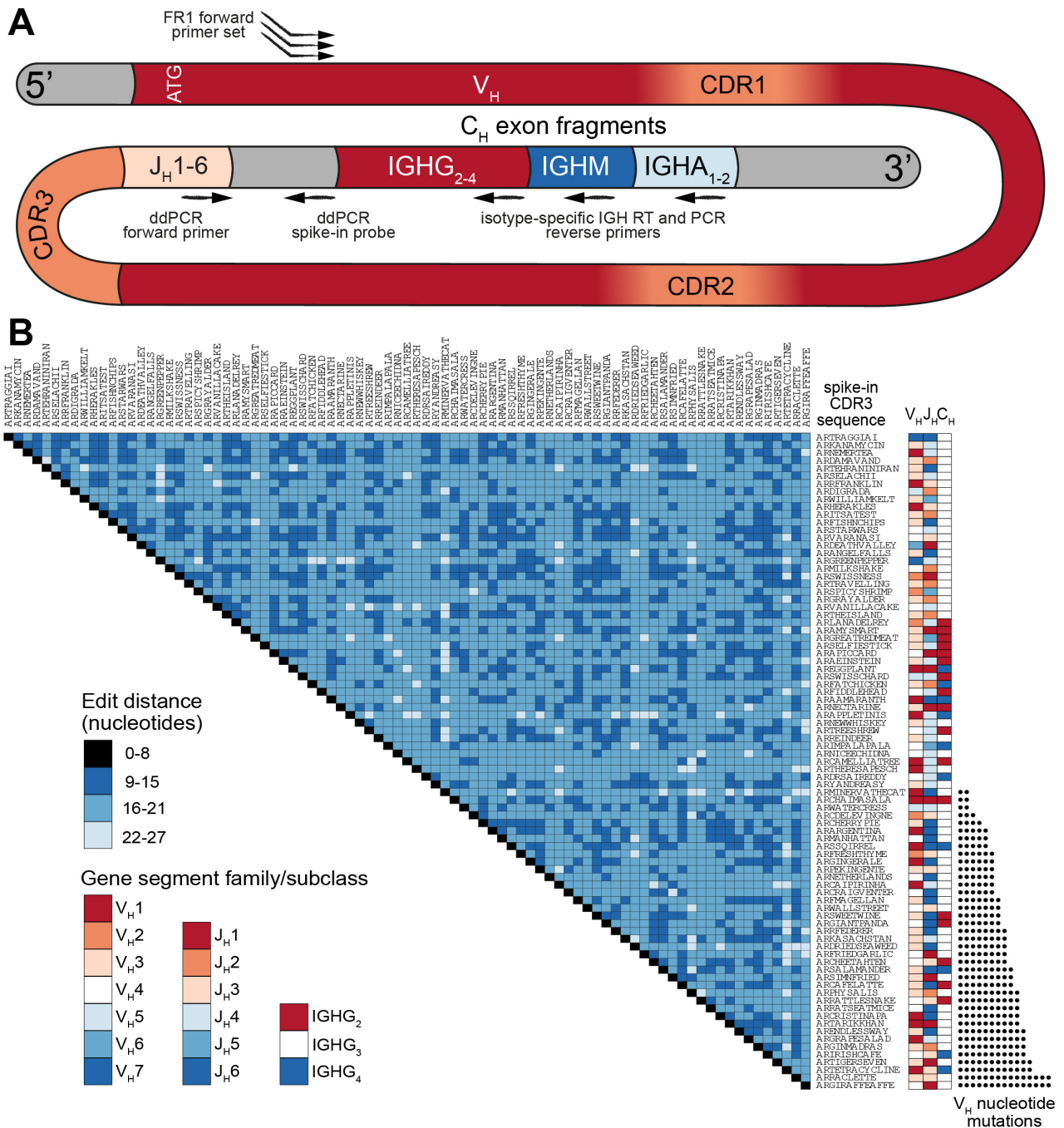


Figure 2. Library preparation of immunoglobulin (Ig) heavy chain genes for high-throughput sequencing (Ig-seq) using molecular amplification fingerprinting (MAF). **(A)** In step 1, reverse transcription (RT) is performed to generate first-strand cDNA with a gene-specific (IgM or IgG) primer which includes a unique reverse molecular identifier (RID) and partial Illumina adapter (IA) region. This results in single-molecule labeling of each cDNA with an RID. In step 2, several cycles of multiplex-PCR are performed using a forward primer set with gene-specific regions targeting heavy chain variable (V_H) framework region 1 (FR1), with overhang regions comprised of a forward unique molecular identifier (FID) and partial IA. In step 3, singleplex-PCR is used to extend the partial IAs. The result (Step 4) is the generation of antibody amplicons with FID, RIDs, and full IA ready for Ig-seq and subsequent MAF-based error and bias correction. **(B)** List of oligonucleotide sequences annealing to the V_H FR1 used in multiplex-PCR (Step 2) of the MAF library preparation protocol. The nearest germline IGHV segment(s) likely to be amplified by the respective primer are listed in the rightmost column. **(C)** The estimated melting temperature distribution of the V_H FR1 forward primer set.

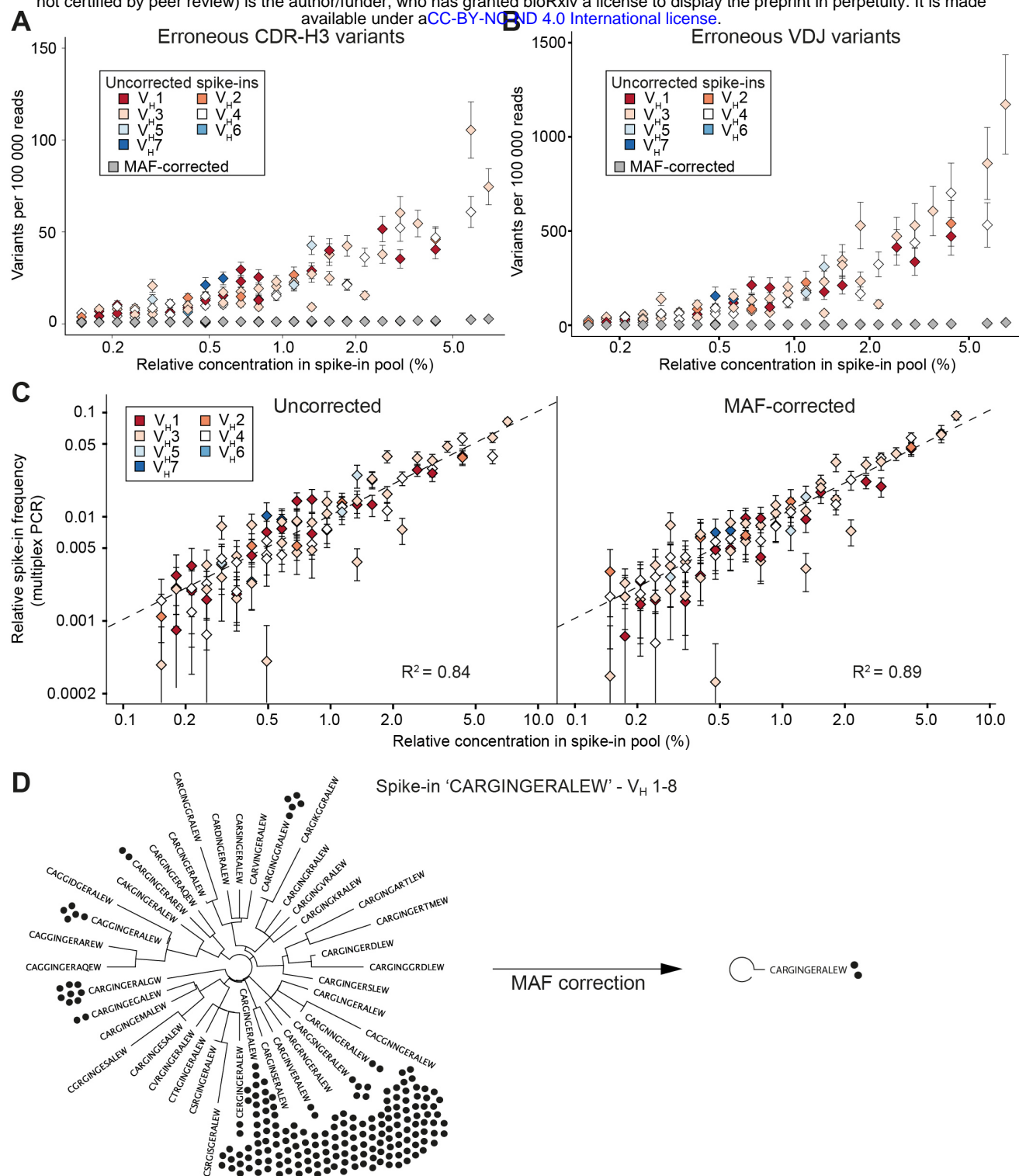


Figure 3. Synthetic standards used to validate performance of error and bias correction of Ig-seq data by MAF. **(A)** The number of erroneous CDR3 sequences (at least one a.a. difference from the correct CDR3) per 100,000 reads is plotted against the relative concentration of each standard (from a master stock). Color-coded diamonds correspond to germline IGHV segment family of the respective standard and show the number of erroneous variants in uncorrected (raw) data; gray diamonds indicate the number of variants remaining after MAF error correction. **(B)** The number of erroneous VDJ variants derived from each standard was calculated by finding all variants that carried the correct CDR3 a.a. sequence, but differed by at least one nt across the entire VDJ region. Colored diamonds represent uncorrected data; gray diamonds indicate variants remaining after MAF error correction. **(C)** Sequencing bias introduced by multiplex-PCR using the FR1 primer set was assessed by plotting the measured frequencies of each standard against its relative concentration (from a master stock). Dashed line represents a bias-free ideal scenario ($R^2 = 1$). The left and right plots show observed frequencies before and after MAF bias correction, respectively. **(D)** Phylogenetic trees visualizing the CDR3 a.a. variants present for a selected standard with the CDR3 a.a. sequence *CARGINGERALEW* and IGHV1-8 and IGHJ1 segment usage. Prior to error correction, 39 erroneous CDR3 a.a. variants (branches) and 218 VDJ nt variants (black circles) were observed. Following MAF error correction, only the original correct CDR3 a.a. and two VDJ nt variants remain. The Ig-seq data sets used in A-C consisted of ~300,000 preprocessed full-length antibody reads from each of the synthetic spike-in only samples. IgG1_D1 dataset was used for panel **(D)** (see **Table S3**).

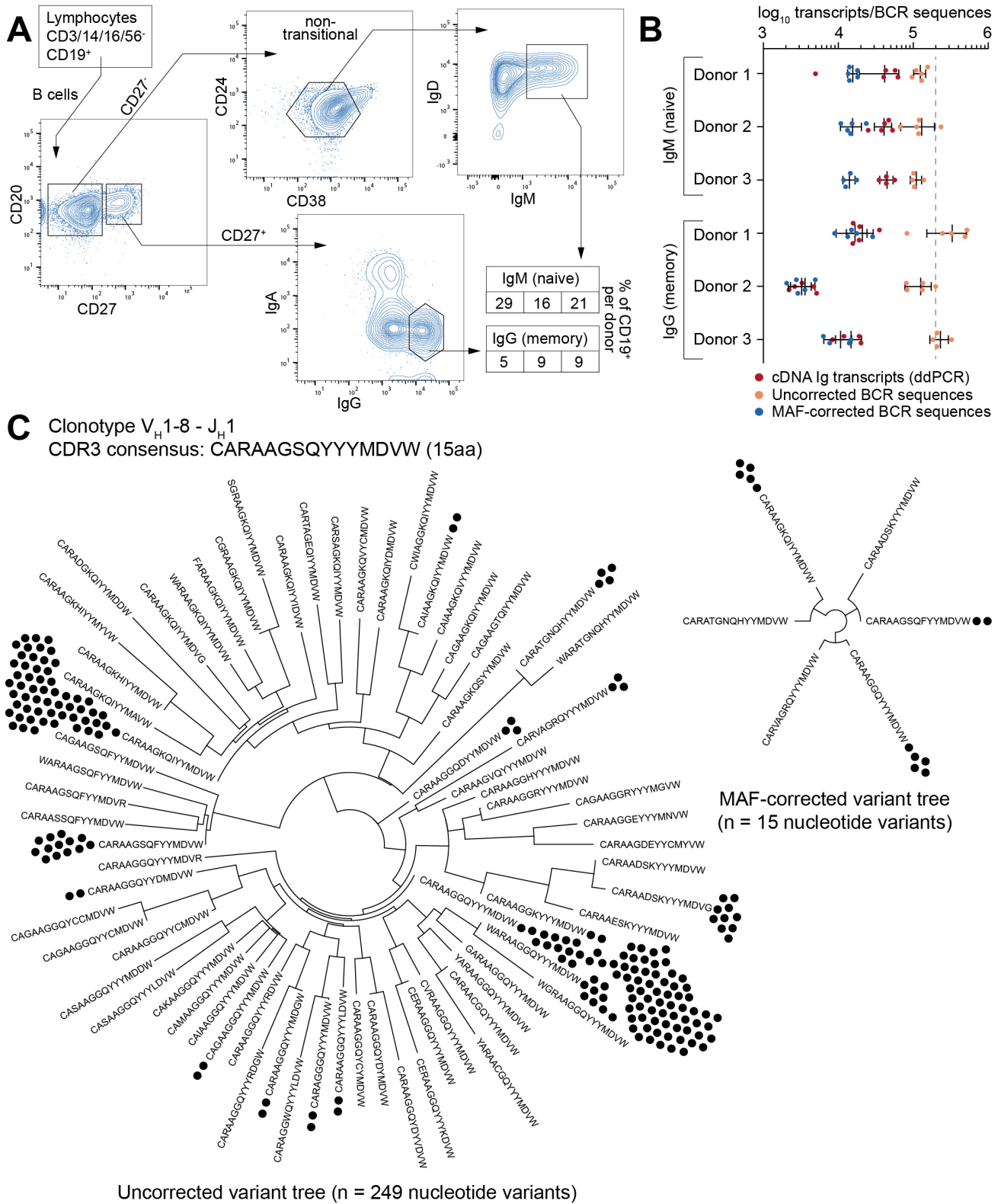


Figure 4. Ig-seq analysis of human naïve (CD27-IgM⁺) and memory (CD27⁺IgG⁺) B cells. **(A)** The flow cytometric workflow for isolating CD27-IgM⁺ and CD27⁺IgG⁺ B cells from peripheral blood. Boxed-in values indicate the frequency of each sorted subset as a percentage of the total B cell (CD19⁺) population from each of three donors (1-3 from left to right). **(B)** Experimental and Ig-seq based quantitation of antibody diversity; points represent cDNA molecule counts (using ddPCR) or unique reads (before and after MAF error correction) from cellular replicates (with mean and standard deviation shown) isolated from each donor and B cell subset. Unique read counts were based on the VDJ nt sequence. Dashed line represents the number of B cells isolated per cellular replicate (2 × 10⁵ cells). **(C)** Phylogenetic trees illustrating CDR3 a.a. and nt variants present for the selected clonotype with the consensus CDR3 sequence *CARAAGSQYYMDVW* and IGHV1-8 to IGHJ1 recombination. Prior to error correction, 70 erroneous CDR3 a.a. variants and 249 VDJ nt variants (black circles) were observed. Following MAF error correction, only 6 CDR3 a.a. and 15 VDJ nt variants remain. The Ig-seq data sets used in **(B)** are described in **Table S3**; IgG1_D1 was used for the tree in panel **(C)**.

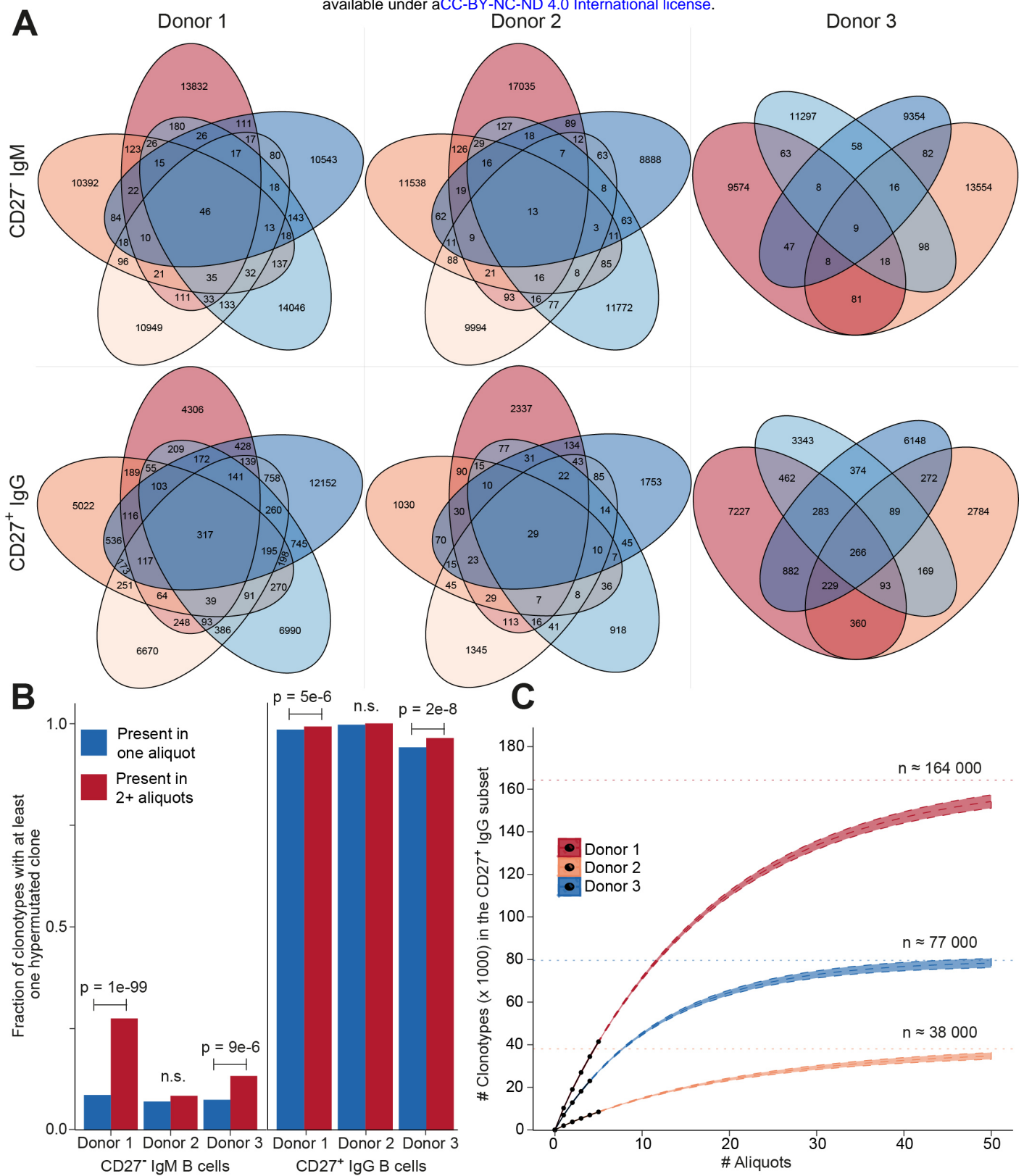


Figure 5. Ig-seq analysis of molecular features highlight global differences between naïve (CD27-IgM⁺) and memory (CD27⁺IgG⁺) B cells. **(A)** Clonotype size (calculated as the total number of variants within a clonotype) for naïve IgM (blue lines) and memory IgG (red lines) B cells. Each pair of lines represents a single donor. **(B)** Graph showing the distribution of average SHM frequencies for VDJ nt variants per clonotype. The CD27-IgM⁺ B cell repertoires (blue lines) have a median SHM value of zero, whereas only a small fraction of clonotypes (approximately 8%) contain an average of one or more mutations. CD27⁺IgG⁺ repertoires (red lines) have a median of 20 to 24 SHM per nt variant within each clonotype. **(C)** CDR3 a.a. length distribution across clonotypes from naïve (blue lines) and memory (red lines) B cell subsets.

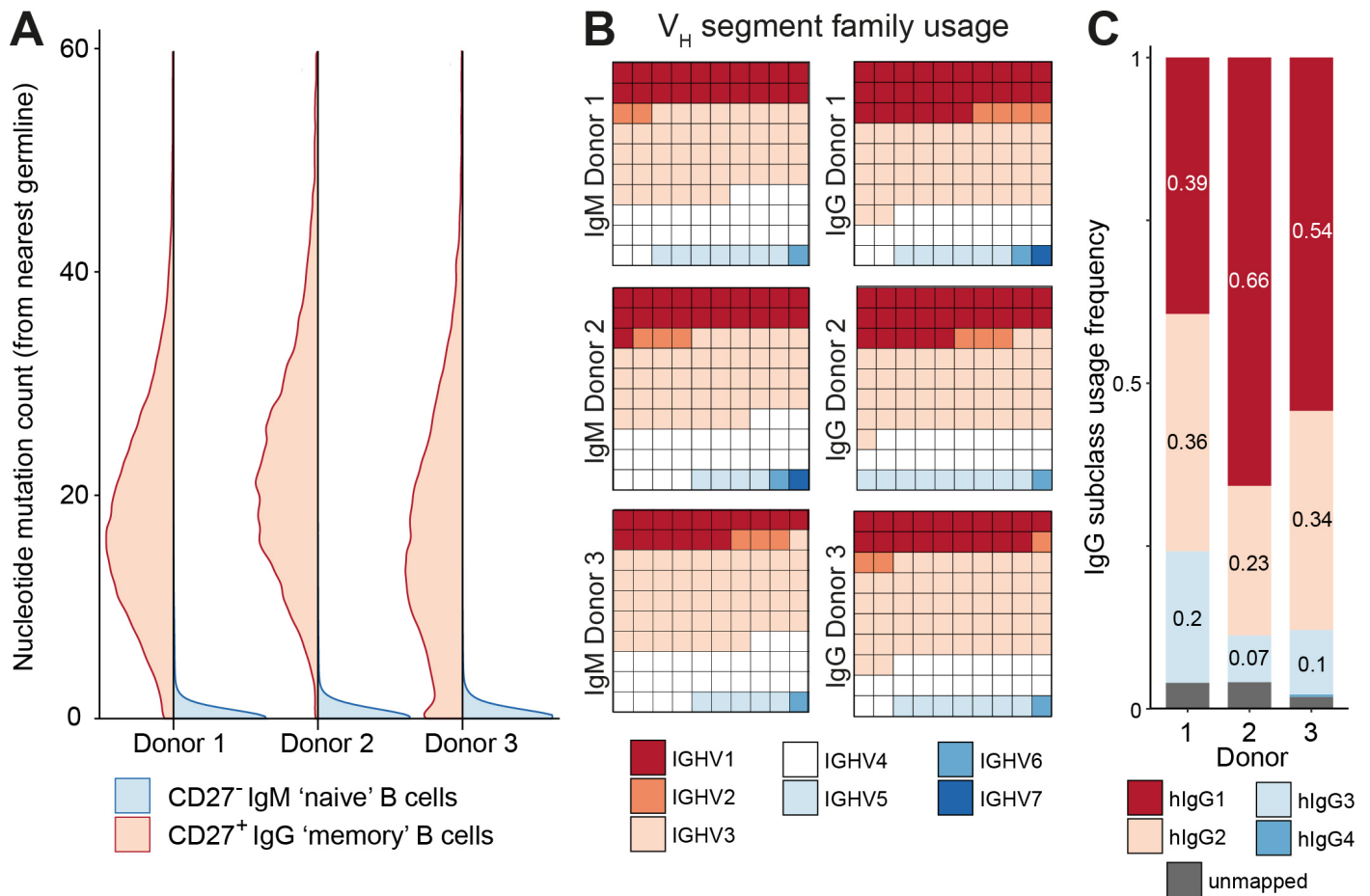


Figure 6. Clonotype diversity analysis across cellular replicates of naïve (CD27⁻IgM⁺) and memory (CD27⁺IgG⁺) B cells. **(A)** Venn diagrams show the presence of clonotypes (80% CDR3 a.a. similarity to least one clone in the cluster, same CDR3 a.a. length, same IGHV and IGHJ gene segment usage) across cellular replicates (2×10^5 cells each) from each donor B cell subset. **(B)** Bar graph showing the fraction of clonotypes containing at least one variant with at least one nt SHM. Blue and red bars indicate clonotypes identified either in only one aliquot (unique) or in several aliquots (shared), respectively. The p-values represent significance using Fisher's exact test. **(C)** Species accumulation curves for CD27⁺IgG⁺ B cells: the number of newly discovered clonotypes from each additional cellular replicate (black circles) is plotted. Extrapolating the observed overlap provides an estimate for the total number of distinct clonotypes (Chao2 estimator: $D1 = 164,268 \pm 2,365$; $D2 = 38,034 \pm 1,302$; $D3 = 76,904 \pm 1,409$) and the approximate amount of cellular replicates needed to discover all clonotypes present in the peripheral blood CD27⁺ IgG⁺ population.

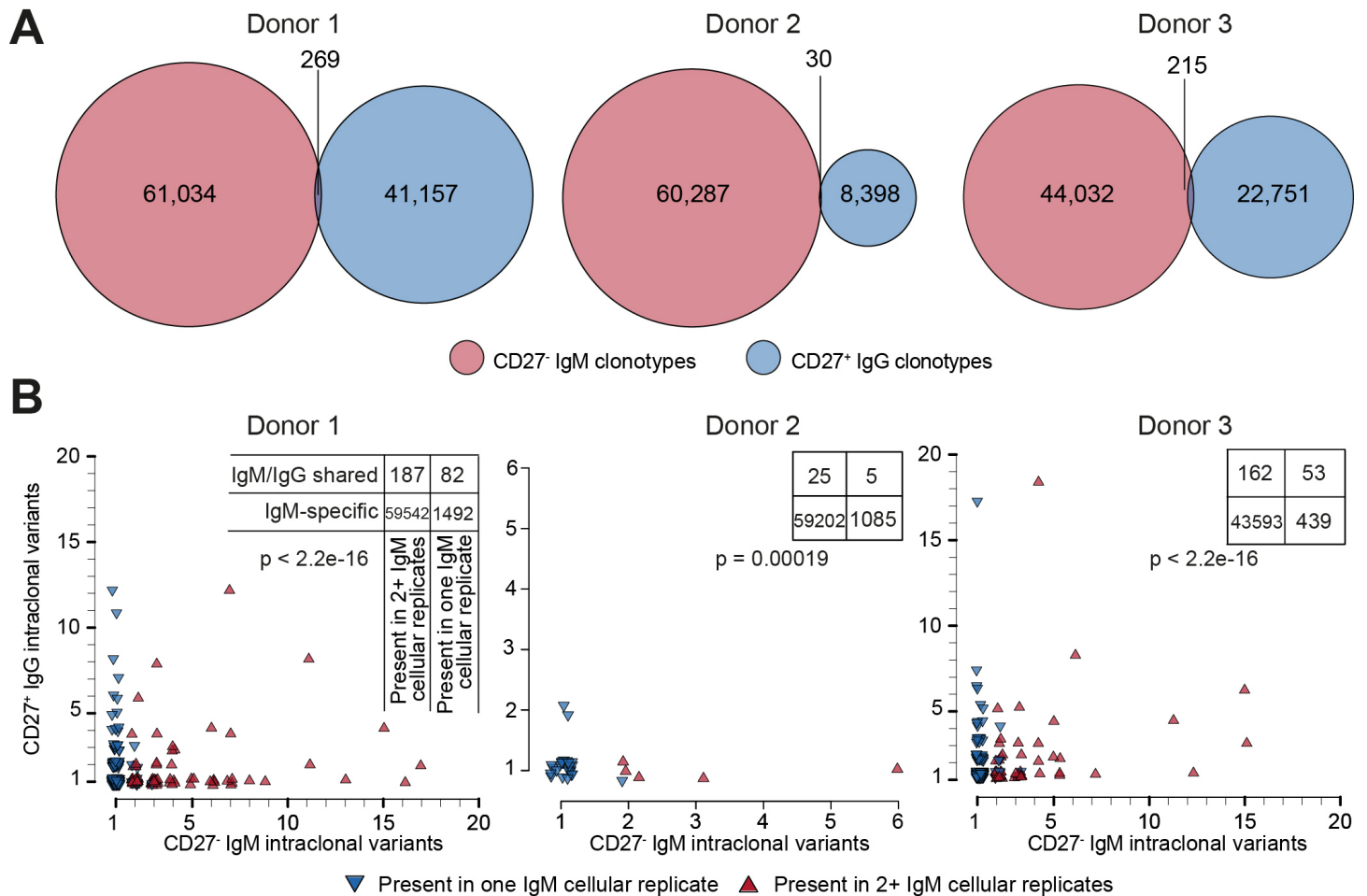


Figure 7. Genetic features of naïve and IgG memory BCR repertoires. **(A)** Block map shows IGHV gene segment usage sorted by family (color-coded) across donors and B cell subset; blocks are normalized to the total number of clonotypes within each group. **(B)** IgG subclass usage in CD27⁺IgG⁺ donor repertoires. IgG₄ sequences (dark blue bars) were virtually absent among three donors; a small fraction of sequences could not be unambiguously mapped based on the sequencing read (gray bars). **(C)** Venn diagrams showing the overlap of clonotypes (80% CDR-H3 amino acid similarity to least one clone in the cluster, same CDR3 a.a. length, same IGHV and IGHJ gene segment usage) between the naïve (CD27-IgM⁺) and memory (CD27⁺IgG⁺) BCR heavy chain repertoire in each of three donors. **(D)** Each plot shows the clonal composition of each shared clonotype (from panel **(C)**) in terms of its IgG and IgM intraclonal variants. Red triangles indicate clonotypes found in multiple IgM cellular aliquots; blue triangles show clonotypes which could only be found in one IgM cellular aliquot. The total number of clonotypes found are depicted in the corresponding contingency table. Fisher's exact test was used to quantitatively analyze enrichment of expanded IgM clonotypes in the shared IgM/IgG subset. **(E)** Ternary plots comprised of three axes representing the IgG1, IgG2, and IgG3 isotype subclasses. The relative subclass composition of intraclonal variants per IgG clonotype (each represented by a circle colored according to the number of variants belonging to that clonotype) is depicted by the position of the circle within the triangular space. Red circles represent the average subclass composition for all IgG variants of each donor (cf. **Figure 6C**).

CDR3 amino acid sequence	CDR3 length	Variable (V) gene segment	Joining (J) gene segment	VH family	JH	IgG subclass	total mismatches	FR1 mismatches	CDR1 mismatches	FR2 mismatches	CDR2 mismatches	FR3 mismatches	spike-in pool contribution (%)
ARYANDREASY	11	IGHV3-23*01	IGHJ4*01	VH3	JH4	IgG3	0	0	0	0	0	0	7.00
ARCAFELATTE	11	IGHV3-23*01	IGHJ6*01	VH3	JH6	IgG2	9	0	4	0	5	0	0.15
ARSIMNFRIED	11	IGHV3-23*01	IGHJ1*01	VH3	JH1	IgG3	9	0	4	1	3	1	1.11
ARDRSAIREDDY	12	IGHV4-34*01	IGHJ4*01	VH4	JH4	IgG4	0	0	0	0	0	0	0.21
ARCRAIGVENTER	13	IGHV4-34*12	IGHJ6*01	VH4	JH6	IgG3	6	0	2	0	3	1	5.92
ARCHEETAHTEN	12	IGHV4-34*12	IGHJ3*01	VH4	JH3	IgG2	8	0	4	0	4	0	0.57
ARTHERESAPESCH	14	IGHV1-69*06	IGHJ4*01	VH1	JH4	IgG3	0	0	0	0	0	0	0.18
ARRSQIRREL	10	IGHV1-69*06	IGHJ6*01	VH1	JH6	IgG4	5	0	1	1	3	0	0.80
ARTARIK KHAN	11	IGHV1-69*06	IGHJ1*01	VH1	JH1	IgG3	10	0	4	0	5	1	3.04
ARCAMELLIATREE	14	IGHV1-18*01	IGHJ4*01	VH1	JH4	IgG2	0	0	0	0	0	0	0.25
ARCISTINAPA	12	IGHV1-18*01	IGHJ6*01	VH1	JH6	IgG3	10	0	4	3	2	1	0.68
ARTETRACYCLINE	14	IGHV1-18*01	IGHJ3*01	VH1	JH3	IgG4	12	0	4	4	4	0	2.57
ARNICEECHIDNA	13	IGHV4-61*08	IGHJ4*01	VH4	JH4	IgG3	0	0	0	0	0	0	0.21
ARGIANTPANDA	12	IGHV4-61*08	IGHJ6*01	VH4	JH6	IgG2	7	0	3	0	4	0	0.94
ARGIRAFFEAFFE	13	IGHV4-61*08	IGHJ1*01	VH4	JH1	IgG3	15	0	4	3	8	0	3.04
ARIMPALAPALA	12	IGHV4-59*01	IGHJ5*01	VH4	JH5	IgG4	0	0	0	0	0	0	0.25
ARRATSEATMICE	13	IGHV4-59*01	IGHJ6*01	VH4	JH6	IgG3	10	0	4	0	4	2	0.48
ARRATTLESNAKE	13	IGHV4-59*01	IGHJ3*01	VH4	JH3	IgG2	10	1	5	0	3	1	2.17
ARREINDEER	10	IGHV3-30*03	IGHJ4*01	VH3	JH4	IgG3	0	0	0	0	0	0	0.18
ARSALAMANDER	12	IGHV3-30*03	IGHJ6*01	VH3	JH6	IgG4	9	0	3	2	4	0	1.32
ARTIGERSEVEN	12	IGHV3-30*03	IGHJ1*01	VH3	JH1	IgG3	12	0	4	0	5	3	2.17
ARTREESHREW	11	IGHV4-39*01	IGHJ4*01	VH4	JH4	IgG2	0	0	0	0	0	0	0.29
ARMANHATTAN	11	IGHV4-39*01	IGHJ6*01	VH4	JH6	IgG3	5	0	1	0	4	0	0.94
ARIRISHCAFE	11	IGHV4-39*01	IGHJ3*01	VH4	JH3	IgG4	11	2	3	2	4	0	1.84
ARNEWWHISKEY	12	IGHV3-48*02	IGHJ4*01	VH3	JH4	IgG3	0	0	0	0	0	0	0.35
ARSWETWINE	11	IGHV3-48*02	IGHJ6*01	VH3	JH6	IgG2	7	0	3	0	4	0	1.56
ARGINMADRAS	11	IGHV3-48*02	IGHJ2*01	VH3	JH2	IgG3	11	0	4	1	6	0	1.84
ARAPPLETINIS	12	IGHV1-3*02	IGHJ4*01	VH1	JH4	IgG4	0	0	0	0	0	0	0.35
ARCAPIRINHA	12	IGHV1-3*02	IGHJ4*01	VH1	JH4	IgG3	6	0	2	0	4	0	1.32
ARNECTARINE	11	IGHV3-21*01	IGHJ1*01	VH3	JH1	IgG2	0	0	0	0	0	0	0.41
ARPHYSALIS	10	IGHV3-21*01	IGHJ2*01	VH3	JH2	IgG3	10	0	5	0	3	2	0.80
ARAAMARANTH	11	IGHV1-2*02	IGHJ6*01	VH1	JH6	IgG4	0	0	0	0	0	0	0.57
ARGRAPESALAD	12	IGHV1-2*02	IGHJ4*01	VH1	JH4	IgG3	11	0	7	0	3	1	0.21
ARFIDDLEHEAD	12	IGHV4-31*02	IGHJ4*01	VH4	JH4	IgG2	0	0	0	0	0	0	0.48
ARFRIEDGARLIC	13	IGHV4-31*02	IGHJ1*01	VH4	JH1	IgG3	8	0	2	0	6	0	0.15
ARFATCHICKEN	12	IGHV3-33*01	IGHJ2*01	VH3	JH2	IgG4	0	0	0	0	0	0	0.48
ARDRIEDSEAWEEED	14	IGHV3-33*01	IGHJ6*01	VH3	JH6	IgG3	8	0	3	0	1	4	1.56
ARSWISSCHARD	12	IGHV5-51*01	IGHJ4*01	VH5	JH4	IgG2	0	0	0	0	0	0	0.29
ARWATERCRESS	12	IGHV5-51*01	IGHJ4*01	VH5	JH4	IgG3	2	0	2	0	0	0	1.11
AREGGPLANT	10	IGHV1-46*01	IGHJ1*01	VH1	JH1	IgG4	0	0	0	0	0	0	0.41
ARARGENTINA	11	IGHV1-46*01	IGHJ6*01	VH1	JH6	IgG3	5	0	3	0	1	1	0.18
ARAEINSTEIN	11	IGHV3-7*01	IGHJ4*01	VH3	JH4	IgG2	0	0	0	0	0	0	0.29
ARKASACHSTAN	12	IGHV3-7*01	IGHJ4*01	VH3	JH4	IgG3	8	0	4	0	4	0	0.94
ARAPICCARD	10	IGHV4-38-2*02	IGHJ1*01	VH4	JH1	IgG2	0	0	0	0	0	0	0.25
ARNETHERLANDS	13	IGHV4-38-2*02	IGHJ6*01	VH4	JH6	IgG3	6	0	1	1	4	0	0.80
ARSELFISTICK	13	IGHV3-11*01	IGHJ4*01	VH3	JH4	IgG2	0	0	0	0	0	0	0.41
ARPEKINGENTE	12	IGHV3-11*01	IGHJ3*01	VH3	JH3	IgG3	6	0	3	1	2	0	1.32
ARCHAIMASALA	12	IGHV1-8*01	IGHJ1*01	VH1	JH1	IgG2	2	0	0	2	0	0	0.21
ARINGERALE	11	IGHV1-8*01	IGHJ6*01	VH1	JH6	IgG3	6	1	0	4	1	0	0.68
ARGREATREDMEAT	14	IGHV3-66*03	IGHJ5*01	VH3	JH5	IgG2	0	0	0	0	0	0	0.68
ARRACLETTE	10	IGHV3-66*03	IGHJ3*01	VH3	JH3	IgG3	15	0	3	2	6	4	0.25
ARAMYSMART	10	IGHV3-30*04	IGHJ1*01	VH3	JH1	IgG2	0	0	0	0	0	0	0.29
ARENDESSWAY	12	IGHV3-30*04	IGHJ6*01	VH3	JH6	IgG3	11	1	4	4	2	0	0.80
ARLANADELREY	12	IGHV2-5*01	IGHJ4*01	VH2	JH4	IgG2	0	0	0	0	0	0	4.24
ARCDELEIVINGNE	13	IGHV2-5*01	IGHJ3*01	VH2	JH3	IgG3	3	0	3	0	0	0	1.11
ARTHEISLAND	11	IGHV3-64*02	IGHJ2*01	VH3	JH2	IgG3	0	0	0	0	0	0	3.59
ARRFEDERER	10	IGHV3-64*02	IGHJ6*01	VH3	JH6	IgG3	8	0	4	1	3	0	0.57
ARVANILLACAKE	13	IGHV4-4*07	IGHJ4*01	VH4	JH4	IgG3	0	0	0	0	0	0	0.35
ARWALLSTREET	12	IGHV4-4*07	IGHJ3*01	VH4	JH3	IgG3	7	0	2	0	5	0	1.11
ARGRAYALDER	11	IGHV3-15*01	IGHJ2*01	VH3	JH2	IgG3	0	0	0	0	0	0	1.84
ARCHERRYPIE	11	IGHV3-15*01	IGHJ6*01	VH3	JH6	IgG3	4	0	0	2	2	0	0.25
ARSPICYSHRIMP	13	IGHV2-26*01	IGHJ5*01	VH2	JH5	IgG3	0	0	0	0	0	0	0.68
ARFRESHTHYME	12	IGHV2-26*01	IGHJ3*01	VH2	JH3	IgG3	6	0	4	0	2	0	0.41
ARTRAVELLING	12	IGHV3-74*01	IGHJ2*01	VH3	JH2	IgG3	0	0	0	0	0	0	0.94
ARFMAGELLAN	11	IGHV3-74*01	IGHJ6*01	VH3	JH6	IgG3	7	0	1	2	4	0	0.35
ARSWISSNESS	11	IGHV2-70*13	IGHJ1*01	VH2	JH1	IgG3	0	0	0	0	0	0	0.15
ARMILKSHAKE	11	IGHV4-30-4*01	IGHJ2*01	VH4	JH2	IgG3	0	0	0	0	0	0	0.48
ARGREENPEPPER	13	IGHV7-4-1*02	IGHJ4*01	VH7	JH4	IgG3	0	0	0	0	0	0	0.57
ARANGELFALLS	12	IGHV3-53*01	IGHJ6*01	VH3	JH6	IgG3	0	0	0	0	0	0	0.29
ARDEATHVALLEY	13	IGHV6-1*01	IGHJ1*01	VH6	JH1	IgG3	0	0	0	0	0	0	0.41
ARVARANASI	10	IGHV4-30-2*03	IGHJ3*02	VH4	JH3	IgG3	0	0	0	0	0	0	0.35
ARSTARWARS	10	IGHV4-28*01	IGHJ4*01	VH4	JH4	IgG3	0	0	0	0	0	0	4.24
ARFISHNCHIPS	12	IGHV3-9*01	IGHJ6*01	VH3	JH6	IgG3	0	0	0	0	0	0	5.92
ARITSATEST	10	IGHV3-20*01	IGHJ2*01	VH3	JH2	IgG3	0	0	0	0	0	0	2.57
ARHERAKLES	10	IGHV1-24*01	IGHJ3*01	VH1	JH3	IgG3	0	0	0	0	0	0	4.24
ARWILLIAMKELT	13	IGHV3-49*03	IGHJ5*01	VH3	JH5	IgG3	0	0	0	0	0	0	3.04
ARMINERVATHECAT	15	IGHV1-69-2*01	IGHJ6*01	VH1	JH6	IgG3	2	0	0	0	2	0	1.56
ARTARDIGRADA	9	IGHV5-10-1*02	IGHJ2*01	VH5	JH2	IgG3	0	0	0	0	0	0	1.32
ARRFRANKLIN	11	IGHV1-58*02	IGHJ3*01	VH1	JH3	IgG3	0	0	0	0	0	0	0.80
ARSELACHII	10	IGHV3-72*01	IGHJ4*01	VH3	JH4	IgG3	0	0	0	0	0	0	0.57
ARTEHRANINIRAN	14	IGHV3-73*02	IGHJ6*01	VH3	JH6	IgG3	0	0	0	0	0	0	0.18
ARDAMAVAND	10	IGHV3-13*01	IGHJ2*01	VH3	JH2	IgG3	0	0	0	0	0	0	0.21
ARNEMERTEA	10	IGHV1-45*02	IGHJ4*01	VH1	JH4	IgG3	0	0	0	0	0	0	0.48
ARKANAMYGIN	11	IGHV3-43*01	IGHJ4*01	VH3	JH4	IgG3	0	0	0	0	0	0	0.68
ARTRAGGIAI	10	IGHV7-81*01	IGHJ6*01	VH7	JH6	IgG3	0	0	0	0	0	0	0.48

Supplementary table 1. Molecular characteristics of 85 synthetic human IgH gene standards ('spike-ins') used in this study. See Figure 1A for spike-in construction schematic.

Sample	Total RNA extracted (ng/ul)	Total RNA (total amount)	Average number of cDNA transcripts per ul	Total number of cDNA transcripts	Sequencing depth (number of reads prior to pre-processing)	Number of reads after pre-processing
IgM1_D1	6.6	132	5'109	62'581	676'557	403'271
IgM2_D1	59.1	1182	4'691	57'469	781'452	508'998
IgM3_D1	10.1	202	3'503	42'916	2'236'338	1'039'415
IgM4_D1	30.6	612	405	4'955	792'330	437'560
IgM5_D1	19	380	3'293	40'333	1'519'422	529'933
IgM1_D2	3.7	74	3'390	41'528	1'925'305	1'111'250
IgM2_D2	53.3	1066	4'330	53'043	673'482	379'811
IgM3_D2	34.4	688	2'048	25'088	992'456	488'230
IgM4_D2	2.2	44	3'075	37'669	981'274	452'817
IgM5_D2	47.4	948	3'855	47'220	2'279'624	440'083
IgM1_D3	7.7	154	4'631	56'734	545'947	333'492
IgM2_D3	29.5	590	3'568	43'702	1'448'340	636'432
IgM3_D3	2.6	52	2'860	35'035	649'626	407'914
IgM4_D3	7.2	144	3'567	43'696	986'419	529'860
IgG1_D1	61.2	1224	1'277	15'643	4'282'524	1'005'440
IgG2_D1	8.7	174	1'289	15'784	2'513'370	2'099'114
IgG3_D1	12.1	242	1'549	18'969	720'804	635'650
IgG4_D1	9.4	188	1'577	19'312	1'057'258	919'252
IgG5_D1	30.7	614	2'884	35'325	1'497'219	1'306'234
IgG1_D2	52.4	1048	421	5'157	1'049'571	913'558
IgG2_D2	3.7	74	181	2'211	959'085	700'808
IgG3_D2	8.5	170	268	3'277	910'336	793'379
IgG4_D2	6.5	130	214	2'622	747'740	662'657
IgG5_D2	5.6	112	386	4'722	998'038	826'427
IgG1_D3	7	140	1'552	19'012	604'835	490'991
IgG2_D3	5.4	108	683	8'367	804'122	658'452
IgG3_D3	5.9	118	972	11'901	756'943	649'224
IgG4_D3	6.9	138	1'639	20'072	2'724'771	1'179'968

Sample	Aligned reads (consensus build)	Raw CDR3 variants (AA)	Raw CDR3 variants w.o. singletons (AA)	Raw total variants (whole VDJ, nt)	MAF corrected CDR3 variants (AA)	MAF corrected clonotypes (Same V/J Gene, same)	MAF corrected variants (whole VDJ, nt)	Donor
IgM1_D1	201'889	40'584	19'815	109'404	14'891	14'629	18'120	D1
IgM2_D1	260'867	37'987	14'952	129'770	11'257	11'088	13'629	D1
IgM3_D1	393'067	47'447	17'357	155'913	11'759	11'639	14'106	D1
IgM4_D1	203'127	46'001	21'116	128'936	15'139	14'919	18'001	D1
IgM5_D1	175'215	32'357	14'870	96'251	11'375	11'185	14'088	D1
IgM1_D2	534'227	65'874	24'260	234'442	18'087	17'646	23'645	D2
IgM2_D2	214'469	37'837	16'146	117'594	12'255	12'055	14'567	D2
IgM3_D2	244'601	35'706	14'292	120'009	10'643	10'439	13'289	D2
IgM4_D2	204'029	37'669	16'697	111'390	12'491	12'269	15'254	D2
IgM5_D2	135'195	26'236	12'557	66'889	9'397	9'292	10'797	D2
IgM1_D3	179'806	32'421	12'987	97'529	9'960	9'808	12'518	D3
IgM2_D3	285'749	42'979	18'424	138'162	14'098	13'866	17'610	D3
IgM3_D3	197'051	33'712	15'530	101'684	11'748	11'567	14'455	D3
IgM4_D3	248'200	31'542	12'514	103'246	9'693	9'582	11'799	D3
IgG1_D1	207'903	32'401	9'126	82'538	7'161	6'756	8'963	D1
IgG2_D1	1'484'045	159'576	36'883	520'865	8'549	7'744	13'393	D1
IgG3_D1	489'704	83'795	19'308	246'780	10'986	9'957	15'014	D1
IgG4_D1	712'466	105'073	24'949	322'613	11'320	10'286	17'581	D1
IgG5_D1	1'025'422	162'369	38'540	489'854	18'690	16'575	28'718	D1
IgG1_D2	644'950	56'789	13'373	199'164	3'233	3'011	4'954	D2
IgG2_D2	429'902	24'874	6'368	80'817	1'650	1'456	2'767	D2
IgG3_D2	541'960	38'086	8'957	134'139	1'965	1'849	2'895	D2
IgG4_D2	472'428	25'158	5'943	83'147	1'398	1'288	2'043	D2
IgG5_D2	566'689	40'047	9'394	132'147	2'501	2'323	3'603	D2
IgG1_D3	367'556	74'198	16'725	207'180	10'665	9'808	14'843	D3
IgG2_D3	485'151	55'285	12'829	182'680	4'559	4'266	6'327	D3
IgG3_D3	502'248	60'245	14'155	206'903	5'430	5'082	7'687	D3
IgG4_D3	693'622	96'133	22'466	325'211	9'239	8'556	14'028	D3

Supplementary table 3. Overview over our experimental results

Primer	Sequence	Notes
IgG_1r	TTGGCACCCGAGAATTCCTACTGHHHHHACAHHHHHACAHHHHNATTGTTCTGGGAAGTAGTCCTTGACCAG	Red part indicates primer binding site, black part contains unique identifier and overhang
IgM_1r	TTGGCACCCGAGAATTCCTACTGHHHHHACAHHHHHACAHHHHNATTACGAGGGGGAAAAGGGTTGG	see above
VH1a	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGGCAGCTGGTGCAGTCTGGGG	see above
VH1b	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTCAGCTGGTGCAGTCTGGAG	see above
VH1c	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGCCAGCTTGTGCAGTCTGGGG	see above
VH1d	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGGCAGCTGGTGCAGTCTGGGC	see above
VH2a	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGCCCAGGTCACCTTGAAGGAGTCTG	see above
VH2b	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGCCCAGATCACCTTGAAGGAGTCTG	see above
VH3a	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGTGAGGTGCAGCTGGTGGAGTC	see above
VH3b	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGTGAAGTGCAGCTGGTGGAGTC	see above
VH3c	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGTGAGGTGCAGCTGGTGGAGTC	see above
VH3d	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGTGAGGTGCAGCTGGTGGAGAC	see above
VH4a	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGCAGCTGCAGGAGTCGGG	see above
VH4b	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTACAGCTGCAGGAGTCGGG	see above
VH4c	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGCAGCTGCAGGAGTCCGG	see above
VH4d	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGCAGCTACAGCAGTGGGG	see above
VH5	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGGCAGCTGGTGCAGTCTGGAG	see above
VH6	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTACAGCTGCAGCAGTCAGG	see above
VH7	CGTTCAGAGTTCCTACAGTCCGACGATCHHHHACHHHHACHHHNGCAGTGCAGCTGGTGAATCTGG	see above
PCR1_1r	ACTGGAGTTCCTTGGCACCCGAGAATTCCTACT*G	'_' indicates phosphorothioate bond
PCR2_f	AATGATACGGCGACCACCGAGATCTACACGTTCTACAGTCCGACGAT*C	'_' indicates phosphorothioate bond
PCR2_r	CAAGCAGAAGACGGCATAACGAGATXXXXXXGTGACTGGAGTTCCTTGGCACCC	Red X indicates Illumina index
IgG_probe	/56-FAM/T+C+TT+CCCC+C+TG+G/3IABkFQ/	'+' indicate locked ribonucleic acids, '/56'-FAM = 5' 6-FAM (Fluorescein), '/3IABkFQ/' = 3' Iowa Black FQ
IgM_probe	/56-FAM/C+CCC+AA+CC+C+TTT/3IABkFQ/	'+' indicate locked ribonucleic acids, '/56'-FAM = 5' 6-FAM (Fluorescein), '/3IABkFQ/' = 3' Iowa Black FQ
Spike_probe	/5HEX/CG+T C+T+G ACT +AGA +ACT +C/3IABkFQ/	'+' indicate locked ribonucleic acids, '/56'-FAM = 5' HEX (Hexachlorofluorescein), '/3IABkFQ/' = 3' Iowa Black FQ
ddPCR_f	GGTCACYGTCTCYTCAG	
ddPCR_r	TGGCACCCGAGAATTC	

Supplementary table 3. Overview over all primers and probes used in this study

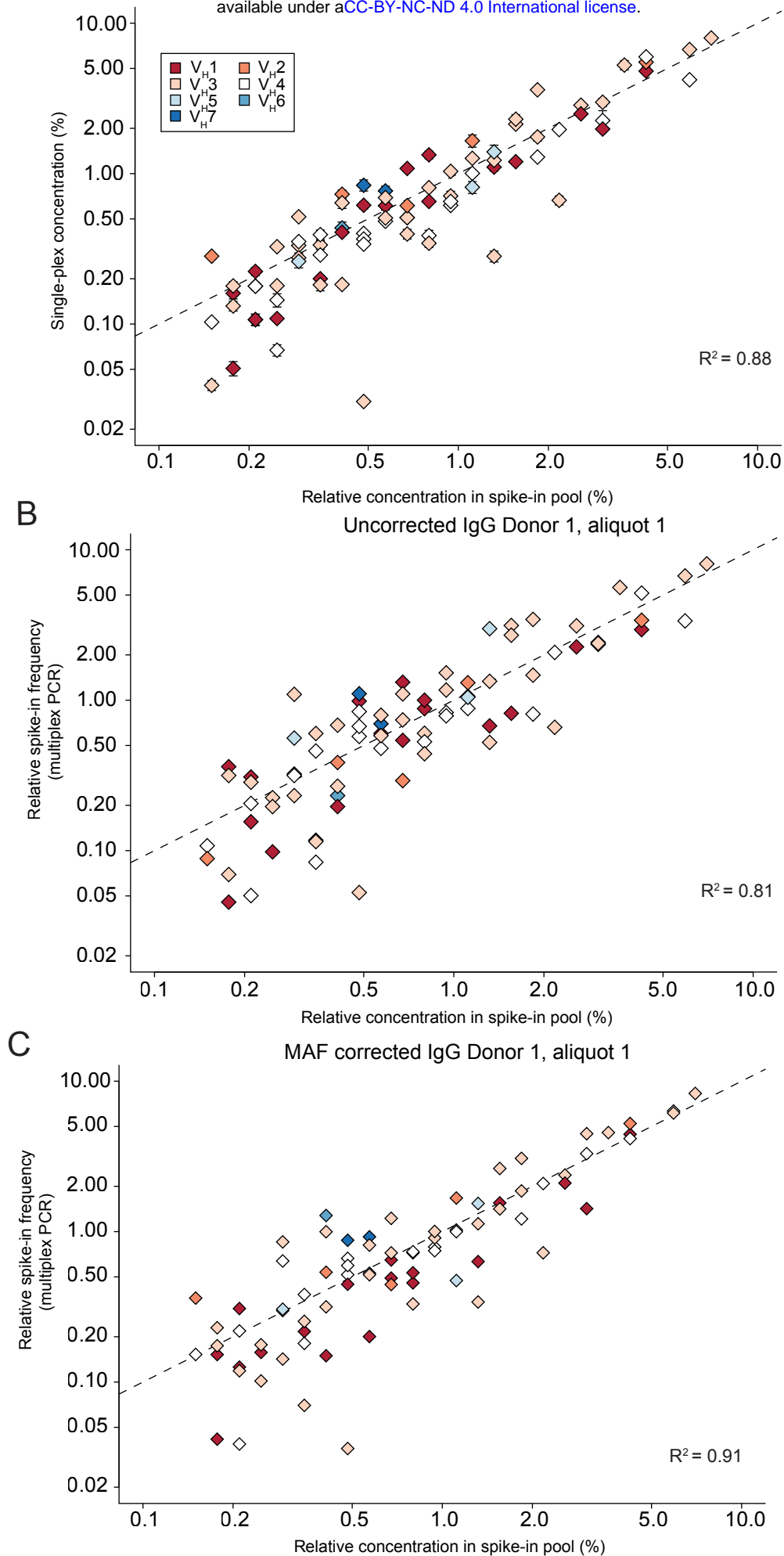
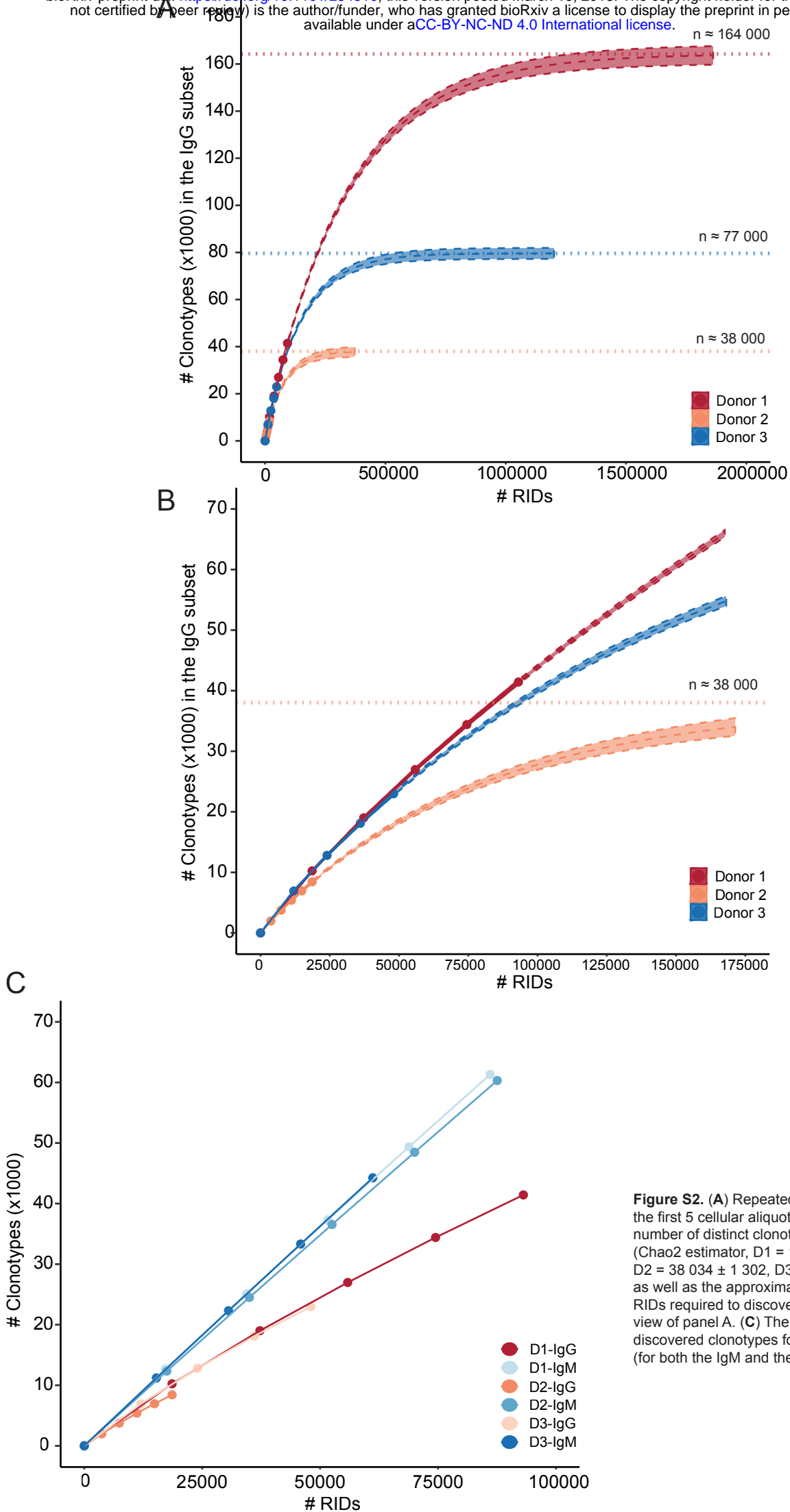


Figure S1. The sequencing bias of the multiplexed primer set was assessed by plotting the measured frequencies of each standard versus its actual, pipetted concentration in the pool. In an ideal case, the measured frequencies would fall onto the dashed line. The deviation from this line was used to calculate the R² value, which decreases with greater deviation. (A) The upper panel shows the divergence of the measured frequencies obtained in the singleplex experiments versus pipetted concentration. Panels (B) and (C) show the measured frequencies and their deviations for one single experiment (IgG, Donor 1, Aliquot 1) before and after MAF correction.



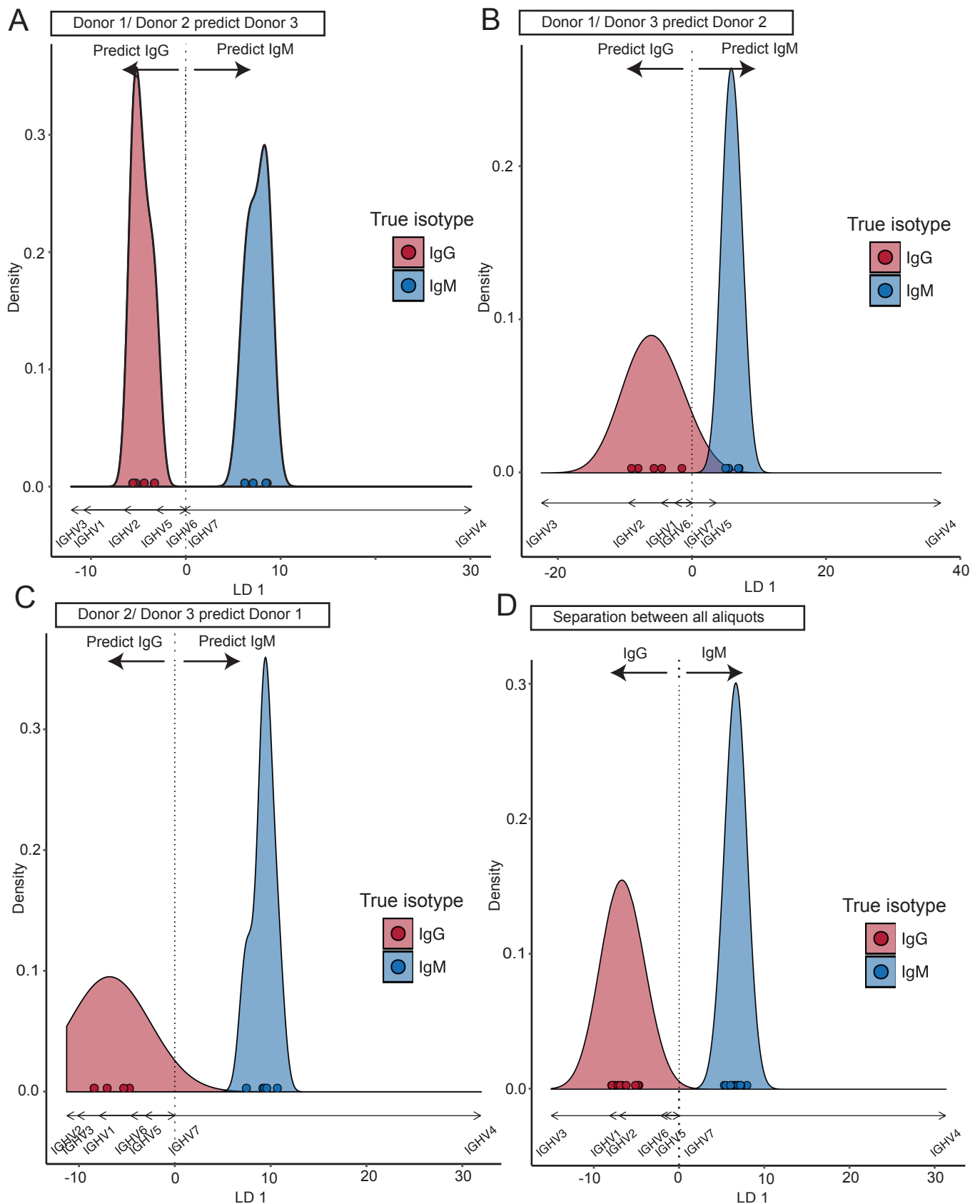


Fig. S3. **Linear discriminant analysis distinguishes CD27+ IgG repertoires and CD27- IgM repertoires based on their V-Gene family usage.** (A) All aliquots from donors 1 and 2 were used to fit an LDA classifier based on the centered log ratio transformed V-Gene family frequencies. Afterwards, the aliquots from donor 3 are projected to the fitted component axis. Positive and negative values predict IgM and IgG repertoires, respectively. Arrows below the plot indicate the contribution of each V-Gene family to the prediction. Colored dots show true class membership and their positions are smoothed using kernel density estimators. Panels (B-D) shows the same procedure for different splits of the data.

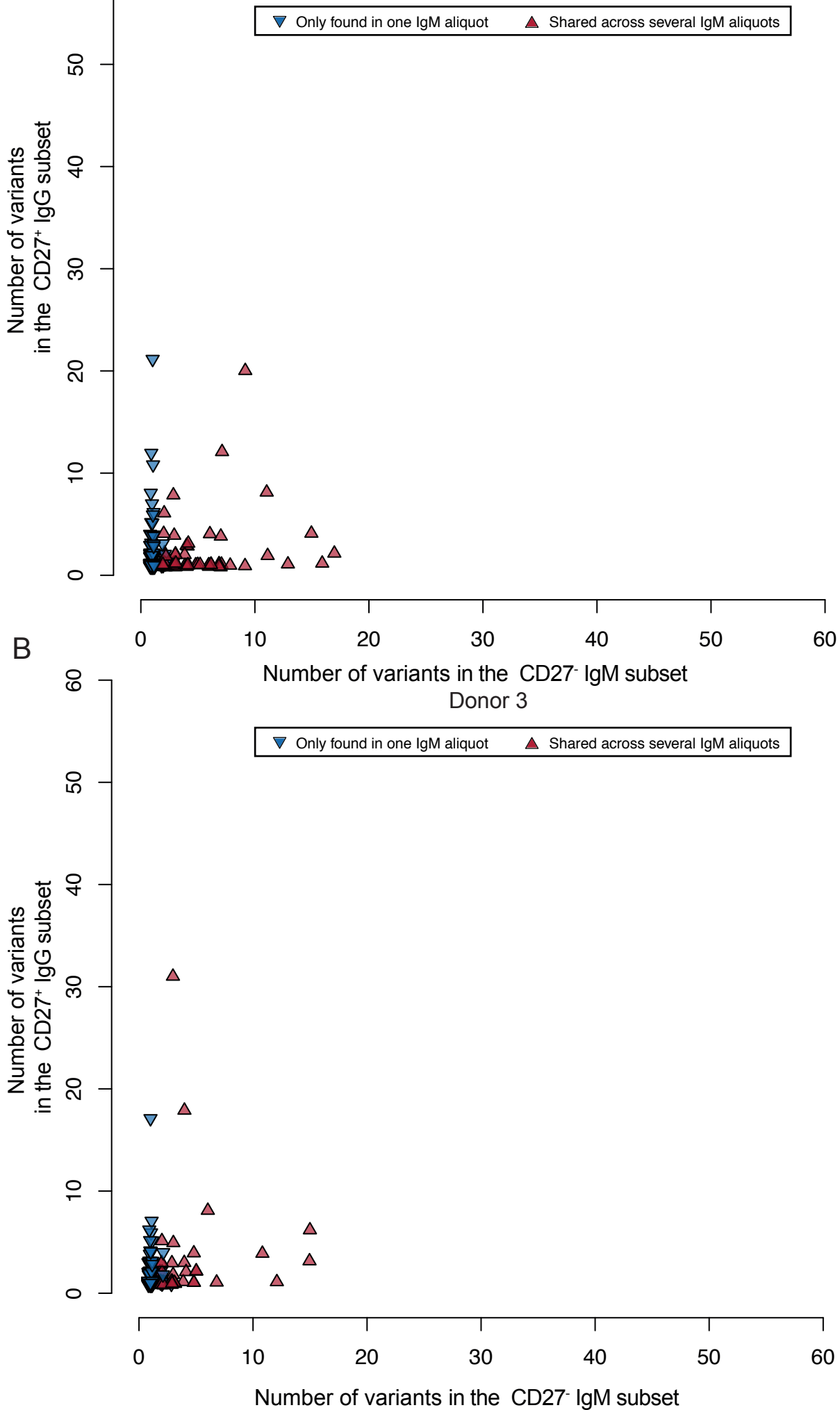


Figure S4. Clonal composition of each clonotype that is shared between the IgG and IgM subsets in terms of its IgG and IgM variants. The red, upward pointing triangle indicates clonotypes that are expanded in the IgM repertoire, whereas the blue triangle highlights clonotypes which could only be found in one IgM aliquot. Panels (A) and (B) show the results for donors 1 and 3, respectively.