



## Network properties derived from deep sequencing of the human B-cell receptor repertoires delineates B-cell populations

Rachael Bashford-Rogers, Anne Palser, Brian Huntly, et al.

*Genome Res.* published online June 6, 2013

Access the most recent version at doi:[10.1101/gr.154815.113](https://doi.org/10.1101/gr.154815.113)

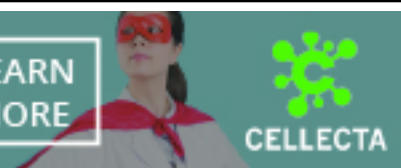
---

<b>P&lt;P</b>	Published online June 6, 2013 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 3.0 Unported), as described at <a href="http://creativecommons.org/licenses/by-nc/3.0/">http://creativecommons.org/licenses/by-nc/3.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---

CRISPR and RNAi Genetic Screening.  
Your new superpower.

LEARN  
MORE



---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

© 2013, Published by Cold Spring Harbor Laboratory Press

1           **NETWORK PROPERTIES DERIVED FROM DEEP SEQUENCING OF THE**  
2           **HUMAN B-CELL RECEPTOR REPERTOIRES DELINEATES B-CELL**  
3           **POPULATIONS.**

4  
5 Rachael J. M. Bashford-Rogers<sup>1</sup>, Anne L. Palser<sup>1</sup>, Brian J. Huntly<sup>2</sup>, Richard Rance<sup>1</sup>, George  
6 S. Vassiliou<sup>1</sup>, George A. Follows<sup>3</sup> and Paul Kellam<sup>1,4</sup>

7  
8 <sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge,  
9 United Kingdom, CB10 1SA;

10 <sup>2</sup>CIMR, Addenbrooke's Hospital, Hills Road, Cambridge, United Kingdom, CB2 0XY;

11 <sup>3</sup>Department of Hematology, Addenbrooke's Hospital, Hills Road, Cambridge, United  
12 Kingdom, United Kingdom, CB2 0QQ;

13 <sup>4</sup>Research Department of Infection, Division of Infection and Immunity, University College  
14 London, Gower Street, London, United Kingdom, WC1E 6BT

15

16 Correspondence should be addressed to Paul Kellam ([pk5@sanger.ac.uk](mailto:pk5@sanger.ac.uk))

17

18 **SCIENTIFIC SECTION DESIGNATION: METHODS AND RESOURCES**

19

20

21 **Abstract**

22           The adaptive immune response selectively expands B- and T-cell clones following  
23 antigen recognition by B- and T-cell receptors (BCR and TCR) respectively. Next-generation  
24 sequencing is a powerful tool for dissecting the BCR and TCR populations at high-resolution,  
25 but robust computational analyses are required to interpret such sequencing. Here, we develop  
26 a novel computational approach for BCR repertoire analysis using established next-generation  
27 sequencing methods coupled with network construction and population analysis. BCR  
28 sequences organize into networks based on sequence diversity, with differences in network  
29 connectivity clearly distinguishing between diverse repertoires of healthy individuals and  
30 clonally expanded repertoires from individuals with chronic lymphocytic leukemia (CLL) and  
31 other clonal blood disorders. Network population measures defined by Gini Index and cluster  
32 sizes quantify the BCR clonality status and are robust to sampling and sequencing depths. We  
33 show, for the first time, B-cell CLL tumor clone evolution over the course of therapy. BCR  
34 network analysis therefore allows the direct and quantifiable comparison of BCR repertoires  
35 between samples and intra-individual population changes between temporal or spatially  
36 separated samples.

37

## 38 **Introduction**

39           Healthy humans have approximately  $3 \times 10^9$  B-cells in the peripheral blood and this  
40 population consists of a repertoire of distinct B-cells expressing different B-cell receptors  
41 (BCRs) necessary to bind diverse antigens and produce an effective humoral immune  
42 response. BCRs consist of two identical heavy chain (IgH) and two identical light-chain  
43 proteins, where the antigen-binding regions are highly diversified (Tonegawa 1983; Woof and  
44 Burton 2004). BCR diversity is generated in a number of ways. The IgH gene locus encodes  
45 for multiple distinct copies of the variable (V), diversity (D) and joining (J) gene segments  
46 (Jung et al. 2006), with functional IgH BCR genes generated by site-specific V-D-J  
47 recombination (Latchman 2005; Schatz and Swanson 2010). The imprecise joining of the V-  
48 D-J gene segments leads to random deletion and insertion of nucleotides during  
49 recombination events, resulting in sequence diversification at the junctional regions (Figure  
50 1A). Rearranged BCR genes are further diversified by helper T-cell mediated somatic  
51 hypermutation (SHM) through the action of activation-induced cytosine deaminase. Through  
52 clonal affinity selection for enhanced antigen binding, non-germline SHM-mediated variation  
53 contributes significantly to the diversification of the mature B-cell repertoire (Brezinschek et  
54 al. 1995; Dorner et al. 1998; Weinstein et al. 2009; Batrak et al. 2011).

55           The diversification and selection dynamics of BCR repertoires in healthy individuals  
56 and those with infection, autoimmunity, immunodeficiency or B-cell malignancies remain  
57 poorly understood but can have important clinical implications. For example, the majority of  
58 B-cells in individuals with B-cell malignancies typically express a single dominant clonal  
59 BCR sequence (Arber 2000; Campbell et al. 2008) and continued intra-clonal tumor evolution  
60 by SHM in patients with B-cell lymphomas has been observed (Stamatopoulos et al. 1996;  
61 Bagnara et al. 2006; Volkheimer et al. 2007). Importantly, patients with chronic lymphocytic  
62 leukemia (CLL) with mutated BCR sequences in the tumor clone compared to the germline  
63 have a prognostically inferior survival rate and requirement of early treatment compared to  
64 those with unmutated malignant clones (Caligaris-Cappio and Ghia 2008). The BCR  
65 sequence repertoire of an individual therefore represents a surrogate of their B-cell clonality

66 status in health and disease, with the potential to give new insights into adaptive immune  
67 response as well as providing diagnostic and prognostic power when used clinically.

68 Previous studies have mainly produced descriptive analyses of the BCR populations.  
69 Isoelectric focusing (IEF) spectrotyping methods (Williamson et al. 1973; Rieben et al. 1996;  
70 Satoh et al. 1996) preceded the advent of sequencing technologies (Arnaout et al. 2011) and is  
71 not quantitative. The potential size of the human repertoire is estimated to be  $10^{11}$  unique  
72 BCR sequences therefore deep, high-throughput sequencing is necessary for sampling this  
73 repertoire robustly and to identify different subsets of BCRs (Dimitrov 2010; Benichou et al.  
74 2012). There are several methods for isolation, amplification and sequencing of B-cell  
75 repertoires. Multiplex PCR amplification, using degenerate PCR primers complementary to  
76 germline V and J segments have been designed and validated previously (van Dongen et al.  
77 2003; Lukowsky et al. 2006; Bruggemann et al. 2007; Evans et al. 2007; van Krieken et al.  
78 2007; Vargas et al. 2008), used in numerous biological studies (Sanchez et al. 2003; Campbell  
79 et al. 2008; Boyd et al. 2009; Boyd et al. 2010; Krause et al. 2011; Jager et al. 2012; Lev et al.  
80 2012; Maletzki et al. 2012) and optimized for clinical use (McClure et al. 2006; Harris et al.  
81 2012; Sproul and Goodlad 2012) although the potential for biased PCR amplification remains.  
82 5' rapid amplification of cDNA ends (5' RACE) has also been used (Bertioli 1997; Freeman  
83 et al. 2009; Varadarajan et al. 2011; Warren et al. 2011), but can suffer from low efficiency  
84 and high levels of non-specific amplification, contamination by short fragments from RNA  
85 degradation or incomplete cDNA synthesis. Both methods, employ PCR and therefore have a  
86 risk of systematic over/underrepresentation of immunoglobulin sequences either through  
87 different primer annealing or different amplification efficiencies of the distinct V families  
88 (Sandberg et al. 2005).

89 Previous studies have qualitatively shown diverse IgH repertoires in healthy patients  
90 contrasting with clonal populations in malignancies (Sanchez et al. 2003; Campbell et al.  
91 2008; Boyd et al. 2009; Carulli et al. 2011; Logan et al. 2011; Maletzki et al. 2012), and that  
92 distinct subsets of B-cells within the same individual have distinct repertoires (Wu et al.  
93 2010). To date next-generation sequencing (NGS) of BCRs have primarily focused on

94 classifying the IgHV, D and J recombination frequencies to understand the diversity of the  
95 BCR repertoire (Stewart et al. 1997; Sanchez et al. 2003; Campbell et al. 2008; Boyd et al.  
96 2009; Weinstein et al. 2009; Boyd et al. 2010; Wu et al. 2010; Jager et al. 2012; Lev et al.  
97 2012; Maletzki et al. 2012). However, computational assignment of V-D-J sequences to  
98 reference databases results in many incompletely assigned IgHV, D and J genes even when  
99 the germline alleles are known (Weinstein et al. 2009). This is most likely due to SHM  
100 masking the identity of the germline genes present in the NGS, or the existence of allelic  
101 variation relative to the reference IgH genes. Further, investigation of V-D-J gene usage  
102 frequencies utilizes only part of the BCR sequence diversity with important information about  
103 the V-D-J joining regions and mutational relationships not considered.

104 Here we propose that analysis of the BCR sequence relationships using the full BCR  
105 V-D-J sequence is more informative for human BCR repertoire analysis than V-D-J gene  
106 classification. We show that human BCR repertoire diversity can be interpreted through full  
107 V-D-J genotype diversity using BCR networks, previously shown to be an intuitive way for  
108 understanding B-cell repertoires in zebrafish (Ben-Hamo and Efroni 2011). In such networks,  
109 the lowest level of organization in a population of B-cells, namely independent B-cells, is  
110 represented by sparse networks whereas highly developed (connected) networks most likely  
111 result from clonal expansions of B-cells, arising through antigenic exposure or B-cell  
112 malignancies (Ben-Hamo and Efroni 2011). Using degenerate PCR-based methods we focus  
113 on sequencing RNA populations to maximize analysis of functionally rearranged BCRs rather  
114 than any non-functional first BCR allele defective rearrangements present in the genomic  
115 DNA from B-cell populations, but with the disadvantage that unequal numbers of RNA  
116 molecules per cell have the potential to inflate or deflate detected B-cell populations in the  
117 repertoire. Through sequencing the BCRs from samples with clonally expanded B-cell  
118 populations (peripheral blood from patients with CLL and human lymphoblastoid cell lines  
119 (LCLs)) as well as diverse BCR populations from peripheral blood from healthy individuals,  
120 we show network analysis provides a robust framework to understand vast sequencing  
121 repertoires by sequence relationships that clearly distinguish between B-cell quantitatively

122 using network measures. This framework is complimentary to existing phylogenetic methods,  
123 and we show, for the first time, B-cell tumor clone evolution over the course of therapy.  
124 These methods are robust to sampling and sequencing depths as well as different sequencing  
125 technologies, thereby allowing the direct comparison of multiple tumor samples from the  
126 same and different patients.  
127

128 **Results**129 **Next generation sequencing of IgH variable genes**

130 We amplified by RT-PCR the expressed rearranged IgHV-D-J loci from mRNA from  
131 human B-cell populations using the consensus IgHJ primer and FR1 or FR2 IgHV family  
132 primers (Figure 1A and Table S1) (van Dongen et al. 2003). Peripheral blood (PB) samples  
133 from thirteen healthy individuals, eleven CLL patients, and eight LCLs yielded PCR products  
134 of expected sizes (310-360bp for FR1 and 250-295bp for FR2 primed samples) and were 454  
135 sequenced (Table 1). Samples yielded an average of 42,324 sequencing reads after filtering  
136 for quality and presence of IgH sequence (Table S2). Two additional samples from CLL  
137 patient A (pre and post treatment) were sequenced on the MiSeq platform. We also analyzed  
138 the BCR 454 sequence datasets from *Boyd et al.* (Boyd et al. 2009), which includes three  
139 healthy individuals and five patients with clonal blood disorders (Table S6).

140 The combined per-base error-rate for the RT-PCR and sequencing process for the 454  
141 platform was similar to other studies (Wang et al. 2007; Boyd et al. 2009) ( $1.74 \times 10^{-4}$ , of  
142 which homopolymeric indels and non-homopolymeric errors accounted for 59.7% ( $1.04 \times 10^{-4}$ )  
143 and 40.3% ( $7.04 \times 10^{-5}$ ) of the total error-rate respectively). Similarly the combined per-base  
144 error-rate for MiSeq was  $2.06 \times 10^{-4}$ .

145 To initially assess the clonality of our samples, we determined the percentage of reads  
146 identical to the most abundant BCR sequence in each sample. The percentage of reads  
147 corresponding to the highest expressed BCR sequence in each of the CLL and LCL samples  
148 (range 39.3%-87.8% and 35.2%-78.7% respectively) were significantly higher than that of PB  
149 from healthy individuals (range 0.10% -14.0%) with a p-value  $<0.001$  (Figure 1B). There was  
150 no significant difference in the percentage of identical reads between the LCL and CLL  
151 patient samples (p-value=0.0594). Therefore, we confirm that the healthy individuals  
152 represent diverse BCR populations, whereas the LCL and CLL samples represent more  
153 restricted or clonal BCR populations. Sanger and MiSeq sequencing confirmed that the  
154 dominant clonal sequences from the CLL samples were identical to that from 454 sequencing  
155 (excluding homopolymeric indels) (Figure S2).

156

157 **Validation of sequencing to represent the B-cell populations**

158 To assess the reproducibility of the RT-PCR sequencing method to sample the BCR  
159 repertoire, we compared the number of overlapping sequences from two types of technical  
160 repeats on a range of samples: i) repeating the RT-PCR and re-sequencing (RT-PCR repeats)  
161 and ii) repeating the 454 sequencing from the same RT-PCR product (sequencing repeats).  
162 The percentage of the sequences shared (no more than 1bp different) between sequencing  
163 runs was calculated using all-against-all alignments. This showed over 98% and 30%  
164 reproducibility for LCL and healthy PB samples respectively (Figure 1C), probably due to the  
165 increased probability of resampling more abundant BCR types in LCLs. It is clear the  
166 sequencing overlaps between RT-PCR repeats (Figure 1D, green bars) are not significantly  
167 different to those between sequencing repeats (Figure 1D, purple bars) ( $p$ -value=0.738 by  
168 paired t-test), suggesting that our RT-PCR amplification and sequencing depth is sufficient to  
169 be representative of the major clonal BCR population in the sample.

170

171 **Comparison between independent primer sets suggests limited primer bias.**

172 To assess whether multiplex PCR methods cause significant PCR amplification bias,  
173 we determined the correlation between IgHV gene usages for samples independently  
174 amplified by the FR1 and/or FR2 primer sets. The IgHV gene usage frequency distributions  
175 for both healthy individuals and LCLs resemble those seen by *Arnaout et al.* (Arnaout et al.  
176 2011) (Figure S3A and B). IgHV gene usage frequencies between FR1 and FR2 amplified  
177 samples from eight LCL samples and four healthy individual samples are highly correlated  
178 ( $R$ -value=0.984, Figure S3C) suggesting there is minimal primer amplification bias.  
179 Performing RT-PCR repeats using the FR1 primer set measures the stochastic effect of  
180 resampling the RNA populations and again we found a strong correlation between IgHV gene  
181 usage frequencies between replicates ( $R$ -value=0.972, Figure S3D). Importantly, for both of  
182 these comparisons, we see a strong linear correlation between IgHV gene usage frequencies  
183 when percentages are greater than 5% of the overall population, but BCR sequences at a

184 frequency of 0-5% show some effect of stochastic sampling, irrespective of primer set usage.  
185 Therefore, overall, we do not see significant primer amplification bias under conditions used  
186 here.

187

### 188 **Limitations of V-D-J classification**

189 We classified the V-D-J genes for each read by sequence similarity to germline  
190 sequences from the ImMunoGeneTics database (IMGT) (Lefranc et al. 2009) (Figure S4).  
191 The majority of sequences could be classified to their most closely related reference  
192 sequences for IgHV and IgHJ genes (average of 99.8% and 96.1% respectively). Substantially  
193 fewer IgHD were identifiable (average of 40.5%) due to the shorter sequence length and  
194 potential insertions and deletions within the joining regions between the V-D-J boundaries,  
195 which has been noted in previous studies (Weinstein et al. 2009). Incomplete V-D-J gene  
196 classification may be due to SHM masking the identity of the germline genes present in  
197 individuals and/or the existence of allelic variants of reference IgH (Boyd et al. 2010). We  
198 find no significant difference between the percentage of classified V, D and J genes of our  
199 dataset compared to that of *Boyd et al.* (Boyd et al. 2009) (Figure S4). To overcome  
200 limitations of IgH V-D-J gene classification, we propose that the use of complete V-D-J  
201 sequence information and mutational relationships would be more informative and robust  
202 framework for BCR repertoire analysis and B-cell population structure than simple V-D-J  
203 gene classification.

204

### 205 **BCR sequences naturally organize into networks based on sequence diversity**

206 For each sample, filtered and trimmed 454 or MiSeq sequences were used directly to  
207 generate a sequence network (Figure 2A). Each vertex represents a different sequence, and  
208 the number of identical BCR sequences defines the vertex size. Edges are created between  
209 vertices that differ by one nucleotide. Clusters are groups of interconnected vertices.  
210 Differences in network architectures are clearly seen by comparing B-cell populations from  
211 healthy individuals and LCLs. In LCLs, the majority of 454 sequences fall within a small

212 number of clusters, as these samples are predominantly comprised of a small number of large  
213 B-cell clone types (Figure 2Bi). In contrast, healthy individuals have sparsely connected  
214 networks where most sequences are unique, thus yielding small vertices indicative of high  
215 BCR sequence diversity in the sampled repertoire (Figure 2Bii). From healthy individuals  
216 4.8-32.2% of BCR sequences fall within clusters of 3 or more reads, with the largest cluster  
217 representing 16.7% (4023 reads) of the total population in healthy individual 10. Sequences  
218 within a cluster are most likely related to a single V-D-J BCR progenitor. Alignment of  
219 sequences within the clusters shows the nucleotide differences are distributed along the length  
220 of the 454 sequences (Figure S5 A-C). In all the healthy individual samples, mutations  
221 significantly occur within the CDR regions, known to be hotspots for somatic hypermutation  
222 (Lin et al. 1997), compared to the FWR regions ( $p$ -value=0.000338, Figure S5D). As  
223 expected, the LCLs showed no significant difference between the mutational proportions of  
224 the CDR and FWR regions. Mismatches are not found primarily at the V-D or D-J boundary  
225 suggesting cluster sequences are derived from the same B-cell precursor.

226 The maximum CLL vertex sizes differ between samples (39.2-99.5% of total  
227 sequences) suggesting that large but variable-sized subsets of B-cells express the predominant  
228 BCR sequence(s), surrounded by BCR variants (including total process errors) of the  
229 dominant sequence. Of note, the extent of cluster size diversity is different between CLL  
230 samples, with some displaying extensive clonal enlargement (Figure 2Ci) whereas others  
231 have more limited clonal expansion (Figure 2Cii) and expansion of two dominant clusters  
232 (Figure 2D). Therefore, the magnitude of connectivity of different samples varies between  
233 individual patients with CLL. However, in all cases, the CLL sequence networks are clearly  
234 distinct from the largely sparsely connected age-matched healthy individuals.

235

### 236 **Population measures capture network and sample diversity**

237 We next aimed to quantify BCR population measures to allow comparison and  
238 interpretation of B-cell repertoire dynamics and biology. We investigated several parameters  
239 to describe different aspects of the B-cell populations. Gini Index is an unevenness measure.

240 When applied to the vertex size distribution for a given sample, these measures indicate the  
241 overall clonal nature of a sample, and when applied to the cluster size distributions, these  
242 measures indicate the overall SHM of a sample. The maximum cluster size is the percentage  
243 of reads corresponding to the largest cluster and indicates the degree of clonal expansion of a  
244 sample. To assess the possibility of dual clonal expansions, we include a measure of the  
245 second maximum cluster size as a percentage of reads in a sample.

246 The LCL samples, due to the more restricted repertoires and highly connected  
247 clusters yield high cluster and vertex Gini Indices (averages of 0.94 and 0.80, range 0.91-0.97  
248 and 0.62-0.91 respectively) (Figure 3A) suggesting high unevenness of the size distributions.  
249 By contrast, B-cell networks of healthy individuals occupy a distinct region of Gini Index  
250 vertex and cluster space (averages of 0.21 and 0.05, range 0.10-0.39 and 0.03-0.11  
251 respectively). The CLL samples occupy a spatial range between healthy individuals and LCL  
252 B-cell population extremes with low vertex (between 0.62 and 0.97), and cluster Gini Indices  
253 (between 0.15 and 0.83), indicating B-cell clonal expansions. There is however considerable  
254 variation between the cluster Gini Indices, with CLL patients 1, 10 and 11 having low cluster  
255 Gini Indices, indicative of a highly expanded dominant cluster or dominant clones. Of note,  
256 one healthy individual (healthy individual 10) has a more developed network as defined by an  
257 increase in connectivity and vertex sizes resulting in higher vertex and cluster Gini Indices  
258 (Figure 3A point (a)). This was also verified by independent sequencing using the FR2 primer  
259 set (Figure S6). Further, the highest expressed BCR sequence for healthy individual 10 has  
260 90.6% sequence identity with the closest germline IgHV gene (16 mismatches in 243bp of  
261 alignment) suggesting that this B-cell clone has undergone SHM.

262 We generated networks from the sequences derived from *Boyd et al.* (Boyd et al.  
263 2009) to validate these population measures on independent BCR sequence data. We show  
264 that the clonal populations of the patients with CLL, small lymphocytic lymphoma (SLL)  
265 and/or follicular lymphoma (FL) are distinct from the diverse populations of healthy  
266 individuals (Figures S7A), occupying equivalent regions of the cluster and vertex Gini Index  
267 graphs to samples within this study. Therefore, we conclude that Gini Index population

268 measure robustly separates distinct B-cell populations into different regions based on the  
269 clonal nature of the sample.

270 We then determined whether we could separate monoclonal expansions, biclonal  
271 expansions and diverse B-cell populations using the maximum cluster sizes and second  
272 maximum cluster sizes (Figure 3B). We show that the CLL and LCL samples have maximum  
273 cluster sizes >30% of the total reads compared to maximum cluster sizes of healthy individual  
274 samples of <20%. However, the LCLs and CLLs collectively occupy two distinct regions in  
275 this space. One group exhibits a single dominant clonal sequence, where the remainder of the  
276 clusters are <5% of the total reads (Figure 3B surrounded by the dashed line). The second  
277 group of samples has two dominant clusters above 40% and 20% of the total reads  
278 respectively. CLL patient 5 repertoire comprises two dominant clonal groups each utilizing  
279 different V-D-J genes ([IGHV3-66\*03/IGHD6-19\*01/IGHJ3\*02] and [IGHV6-1\*01/IGHD3-  
280 3\*01/IGHJ4\*02] respectively), where the two clones are unlinked and represented by  
281 completely different BCR sequences (Figure 2D and S8). Limited polyclonal expansions were  
282 observed also in 5/8 of the LCL samples reflecting that EBV transformation of peripheral B-  
283 cells frequently results in polyclonal LCLs. Using the dataset from *Boyd et al.* (Boyd et al.  
284 2009), we show the same phenomenon of polyclonal expansions in a subset of samples  
285 (patients with CLL/SLL and FL/SLL, Figure 2Eiii) where the maximum cluster sizes are  
286 >35% and second maximum cluster sizes are >19% of the total reads (Figure S7B). Therefore  
287 the polyclonal status of the tumor samples can be determined using B-cell network  
288 reconstruction and analysis.

289 An important requirement of this approach is that the network diversity measures  
290 must not be highly dependent on the depth of sequencing (scale invariant) and volume of PB  
291 sample. If a given diversity measure is scale invariant for B-cell networks then the network  
292 diversity measure should be the same regardless of the depth of sampled sequences, i.e. a  
293 subset of 454 sequences should yield the same network diversity measure as the full set of  
294 sequences. We tested all the proposed population measures as a function of sequencing depth  
295 by randomly sampling different proportions of the sequence data for each sample followed by

296 calculation of the corresponding network parameters for both the vertex and cluster size  
297 distributions for the LCL, CLL and healthy samples. All the proposed measures showed little  
298 variation at different sample sizes even when sub-sampling as low as 20% of the original data  
299 size (Figure S9). Below 20%, small deviations in the Gini Index measures are seen. As these  
300 network measures had minimal standard deviation over all sub-sampling ranges, they are  
301 therefore robust parameters for inter-sample comparison.

302

### 303 **Minimal effect of sequencing errors on network properties**

304 We determined whether clusters were likely to be due to the process of somatic  
305 hypermutation or sensitive to or generated through sequencing error of a unique amplified  
306 BCR sequences. For a given BCR sequenced multiple times, such as when multiple B-cells  
307 express identical BCRs, we estimated the expected number of vertices comprising a cluster  
308 that could be due to sequencing error, given our experimentally derived PCR and sequencing  
309 error-rates. We find that all the samples have cluster sizes greater than that expected due to  
310 error alone, even at twice the measured error-rate (Figure S10). Therefore, the connectivity  
311 patterns of networks predominantly reveal differences in clonal expansions of B-cell  
312 populations rather than total sequencing errors. We propose that clusters identified in BCR  
313 networks are therefore derived from B-cells that share a common pro-B-cell progenitor with  
314 rearranged V-D-J that have subsequently expanded and diversified.

315 Directly comparing the Gini Index measures of V-D-J sequences from samples  
316 amplified independently by distinct primer sets (FR1 or FR2 primer sets) showed a strong  
317 positive linear correlation between the two primer sets, with R-values of 0.999 and 0.996  
318 respectively for the vertex and cluster size diversities (Figure S11A). This supports a lack of  
319 PCR or sampling bias or an effect of sequencing errors for independent RT-PCRs with FR1  
320 compared to FR2 primer sets. Further, we find that networks generated to include edge  
321 lengths of up to 5 base changes faithfully retain the network architecture for both the LCL and  
322 healthy individual samples (Figure S11B).

323

**324 BCR repertoire network properties relates to CLL development**

325 To assess the sensitivity of this analysis method, we use the titration experiment from  
326 *Boyd et al.* (Boyd et al. 2009) in which serial 10-fold dilutions of a known clonal CLL PB  
327 sample into normal peripheral blood was performed. We find 90.9% of all reads in the  
328 undiluted sample fall within the leukemic cluster (Figure 3C and D). Using our methods, we  
329 can detect the leukemic clonal sequences at dilutions as low as 1:100,000. When the leukemic  
330 cluster sequences are unknown, detection of expanded clones relies on detecting the  
331 maximum cluster size that is significantly different from that of healthy individuals. We see  
332 significant increases in maximum cluster size above that of the healthy individual in CLL  
333 dilutions of 1:100 or less. Therefore, deep sequencing of BCR repertoires potentially allows  
334 the detection of a clonal lymphoid population in a background of polyclonal cells without  
335 prior knowledge of the leukemic sequence types by comparing to healthy control BCR  
336 clonality (Sayala et al. 2007).

337 We therefore sought to understand the relationship between the BCR population  
338 measures and the CLL clinical information for each patient. Interestingly, there was a strong  
339 correlation between the length of time since CLL diagnosis and the vertex Gini Index (Figure  
340 3E). This suggests longer disease times lead to larger vertices representing larger tumor clonal  
341 populations, in agreement with previous studies (Kelly et al. 2002; Hayes et al. 2010).

342 By BCR deep sequencing and network analysis, we hypothesize that we can follow  
343 the dynamics of dominant clonal populations at multiple time points. To test this, we sampled  
344 a stage B CLL patient (patient A) during the course of therapy. A pre-treatment sample was  
345 taken immediately before entering second cycle of Chlorambucil treatment, and second  
346 sample taken 28 days later. The BCRs were sequenced on the MiSeq platform yielding  
347 40,414 and 36,197 high-quality reads respectively. The most abundant BCR sequences in the  
348 two samples were identical. Network construction shows a single dominant cluster at both  
349 time points, which decreases from 86% to 53% of total reads, reflecting the reduction in WBC  
350 (Figure 4A). This corresponds to the vertex Gini Index reducing from 0.892 to 0.713 and the  
351 cluster Gini Index reducing from 0.244 to 0.138. A composite network was generated of all

352 the sequences from the dominant clusters in both time-points, where the vertex sizes  
353 correspond to frequencies of each BCR at either the pre or post treatment samples (Figure  
354 4B). The proportions of each unique BCR sequence between the pre and post-treatment  
355 samples give a strong linear relationship ( $R$ -value=0.999995, Figure 4C), where the post-  
356 treatment proportional frequencies are 62% of those in the pre-treatment sample, indicating  
357 that all BCR clones within the leukemic cluster are equally affected by this treatment.

358 To understand the evolution of the leukemic cluster over time, a maximum parsimony  
359 tree was fitted (Schliep 2011) encompassing sequences that were observed at least 6 times for  
360 the two samples (Figure 4D). Bootstrapping was performed to evaluate the reproducibility of  
361 the trees, showing strong tree support (>95% certainty for all branches). 122/231 of the BCR  
362 sequences are unique to the pre-treatment sample (e.g. Figure 4, clone 1) compared to 13/231  
363 unique to the post-treatment sample (e.g. Figure 4, clone 2). However, the sequences that are  
364 primarily observed in the post-treatment sample show divergence from the dominant  
365 leukemic clone, suggesting leukemic cluster BCR evolution during treatment (e.g. clone 2). A  
366 similar analysis was performed on an independent Boyd et al. (Boyd et al. 2009) for a patient  
367 with chronic lymphocytic leukemia and small lymphocytic lymphoma samples separated by 3  
368 months (Figure S12) showing distinct clonal diversification patterns.

369

## 370 Discussion

371 Deep-sequencing of B-cell and T-cell repertoires offers the potential for quantitative  
372 understanding of the adaptive immune system in health and disease. Here we use deep-  
373 sequencing of B-cell receptor V-D-J population frequencies and novel analyses of BCR  
374 repertoires at the level of clonal populations. The observation of frequent multiple identical  
375 BCR sequences in tumors and much lower frequency identical BCR sequences in PB from  
376 healthy individuals suggests that we rarely sequence multiple identical RNA molecules from  
377 a single B-cell. Therefore, clusters of related sequences are likely to represent BCRs from  
378 clonal expansions of evolutionarily related B-cells, whereas naïve B-cell populations form  
379 singletons in sparsely connected networks. The effects of RT-PCR or sequencing error and

380 amplification bias on our analysis, often of concern for deep-sequencing, are minimal. We  
381 show a strong linear correlation between the network parameters of samples that have been  
382 RT-PCR amplified using different primer sets to distinct regions of the IgH variable RNA  
383 transcript suggesting that the PCR methods here have limited primer or amplification bias.  
384 We confirm the dominant clonal sequences for the CLL patients by Sanger sequencing, and  
385 show that in all cases the samples have cluster sizes notably greater than that expected due to  
386 the measured total process error-rate. Therefore, these observed V-D-J clusters are likely to  
387 have undergone mutational processes greater than process errors.

388         We define for the first time B-cell V-D-J sequence population measures that describe  
389 the clonality of the sequences and quantify both the effect of B-cell sequence diversification  
390 (cluster size) and clonal proliferation (vertex size) using the Gini Index as an unevenness  
391 measure. The maximum and second maximum cluster size is used to assess dual clonal  
392 expansions. If the B-cell network from limited sequencing is a random sample of the entire  
393 circulating peripheral blood BCR repertoire, then a scale invariant diversity measure should  
394 also capture the predominant structure of the unsampled network. We show that network  
395 structures, combined with these population measures discriminate between B-cell repertoires  
396 of different clonalities in health and disease. These measures are robust to variations in  
397 sequencing and sampling depth and different filtering strategies and are applicable to  
398 independently produced datasets (Boyd et al. 2009). Using different primer sets, sequencing  
399 depths and sequencing technologies, the samples still cluster according to the clonal nature of  
400 the samples, occupying the equivalent distinct regions of Gini Index and maximum/second  
401 maximum graphs. Therefore this analytical strategy is applicable to any BCR deep  
402 sequencing technology.

403         We observed variation between the BCR repertoires in healthy individuals and in  
404 CLL. One healthy individual showed a more developed network, defined by an increase in  
405 connectivity, with corresponding higher Gini Index values and larger maximum cluster sizes  
406 compared to the other healthy individuals. This clone was not germline in sequence and could  
407 be a result of antigen specific memory B-cell expansion or an undiagnosed malignant

408 transformation. Variation in network structures between individual healthy zebrafish BCR  
409 repertoires was also observed in the study by *Ben-Hamo et al.* (Ben-Hamo and Efroni 2011),  
410 where higher connectivity suggested an immune response within the individual. Similarly in  
411 CLL, assessing the maximum and second maximum cluster sizes we identify patients with  
412 more than one BCR clonal expansion where the two dominant clones have different V-D-J  
413 gene usages. This may be due to either the expansion of two distinct malignant B-cell  
414 transformations, or separate antigen-stimulated B-cell clonal expansion unrelated to CLL.  
415 These methods used in time-series may allow the distinction between antigen-driven positive  
416 selection in CDRs compared to malignant-driven expansion (Figure S5). Multiple separate  
417 clonal B-cell populations have been observed in previously published data in a subset of  
418 patients identified by different V and J chain usages (Hsi et al. 2000; Boyd et al. 2009), but  
419 the clinical significance of these findings are not known.

420 Time-dependent evolution of BCR networks may however provide a powerful means  
421 of assessing B-cell tumor evolution and response to therapy as well as the dynamics of a  
422 healthy B-cell repertoire. CLL vertex Gini Index is correlated with the time an individual has  
423 been living with CLL (Figure 3E). This coupled with the observation of *in vivo* evolution of  
424 BCR clones in CLL (Figure 4) suggests BCR sequencing in CLL may provide an additional  
425 prognostic value for the disease. Divergent evolution from a common leukemic ancestor has  
426 previously been observed in CLL, possibly through the accumulation of driver mutations with  
427 selective advantages in growth over other sub-clones (Campbell et al. 2008). Hypermutations  
428 within the IgH locus may also play a driver role in clonal expansions (Ghia and Caligaris-  
429 Cappio 2006). Therefore BCR sequencing and subsequent network and evolutionary analysis  
430 may play an important role in identifying population changes. However, an evolutionary  
431 model for BCR diversity in health and disease, similar to the models used in infectious  
432 disease phylogenetics is needed to fully explore these possibilities. Nevertheless, for the first  
433 time we show the short-term effect of therapy on the B-cell repertoire in CLL, and  
434 demonstrate how networks lend themselves to phylogenetic approaches. These methods are

435 sensitive and informative for characterizing of B-cell populations in health and B-cell

436 malignancies.

437

438

439 **Methods**440 **Samples**

441 Peripheral blood mononuclear cells (PBMCs) were isolated from 10ml of whole  
442 blood from healthy volunteers and CLL patients using Ficoll gradients (GE Healthcare). Total  
443 RNA was isolated using TRIzol® and purified using RNeasy Mini Kit (Qiagen) including on-  
444 column DNase digestion according to manufacturer's instructions. Total RNA was also  
445 isolated from  $1 \times 10^6$  cells from Human lymphoblastoid cell lines (LCLs) from the HapMap  
446 project (Frazer et al. 2007). Research was approved by relevant institutional review boards  
447 and ethics committees (07/MRE05/44)

448

449 **RT-PCR**

450 RT-PCR reagents were purchased from Invitrogen. The FR1 and FR2 primer sets  
451 used (supplied by Sigma Aldrich) are described by Van Dongen *et al.* (van Dongen et al.  
452 2003) and in Table S1. Reverse transcription was performed using 500ng of total RNA mixed  
453 with 1µl JH reverse primer (10µM), 1µl dNTPs (0.25mM), and RNase free water added to  
454 make a total volume of 11µl. This was incubated for 5 minutes at 65°C, and 4µl First strand  
455 buffer, 1µl DTT (0.1M), 1µl RNaseOUT™ Recombinant Ribonuclease Inhibitor and 1µl  
456 SuperScript™ III reverse transcriptase (200units/µl) was added. RT was performed at 50°C  
457 for 60 minutes before heat-inactivation at 70°C for 15 minutes. PCR amplification of cDNA  
458 (5µl of the RT product) was performed with the JH reverse primer and the FR1 or FR2  
459 forward primer set pools (0.25 µM each), using 0.5µl Phusion® High-Fidelity DNA  
460 Polymerase (Finnzymes), 1µl dNTPs (0.25mM), 1µl DTT (0.25mM), per 50µl reaction. The  
461 following PCR program was used: 3 minutes at 94°C, 35 cycles of 30 seconds at 94°C, 30  
462 seconds at 60°C and 1 minute at 72°C, with a final extension cycle of 7 minutes at 72°C on  
463 an MJ Thermocycler.

464

**465 Sequencing methods**

466 454-libraries were prepared using standard Roche-454 Rapid Prep protocols  
467 incorporating 10-base multiplex identifier (MID) tags and sequenced using an FLX Titanium  
468 Genome Sequencer (Roche/454 Life Sciences) or by 250bp paired-ended MiSeq (Illumina).  
469 Raw 454 or MiSeq reads were filtered for base quality (median >32) using the QUASR  
470 program (<http://sourceforge.net/projects/quasr/>). MiSeq forward and reverse reads were  
471 merged together if they contained identical overlapping region of >65bp, or otherwise  
472 discarded. Non-immunoglobulin sequences were removed and only reads with significant  
473 similarity to reference IgHV genes from the IMGT database (Lefranc et al. 2009) using  
474 BLAST (Altschul et al. 1990) were retained ( $1 \times 10^{-10}$  E-value threshold). Primer sequences  
475 were trimmed from the reads, and sequences retained for analysis only if both primer  
476 sequences were identified and if sequence lengths were greater than 255bp or 195bp for FR1  
477 and FR2 primed samples respectively for 454, or both forward and reverse reads greater than  
478 110 bp for MiSeq. FR1 primed PCR samples from CLL patients were also Sanger-sequenced.

479

**480 Per-base error quantification**

481 The same PCR protocol and read quality filtering was used to amplify beta-actin,  
482 beta-globin and GAPDH genes from two healthy individuals (amplicon sizes of 150bp, 340bp  
483 respectively). The sequence representing the majority of the reads for each sample was  
484 classified as the 'true' gene sequence for that individual to account for individual allelic  
485 variation. Any differences between this sequence and the reads were considered to be PCR  
486 and/or sequencing error and classified as homopolymeric indels (occurring in a region of two  
487 or more consecutive identical bases), non-homopolymeric indels, or mismatches.

488

**489 Reference-based V-D-J assignment**

490 BLAST (Altschul et al. 1990) was used to align the 454 sequences against known  
491 BCR sequences from the ImMunoGeneTics (IMGT) database (Lefranc et al. 2009). Due to

492 the difference in length of the different gene families, different BLAST e-value thresholds  
493 were used for the IgHV, IgHD, and IgHJ-genes ( $10^{-70}$ ,  $10^{-3}$  and  $10^{-20}$  respectively).

494

#### 495 **Network assembly and analysis**

496 The network generation algorithm is summarized in Figures 2A and S1. Briefly, each  
497 vertex represents a unique sequence, where the relative size of the vertex is proportional to  
498 the number of sequence reads identical to the vertex sequence. Edges were calculated  
499 between vertices that differed by single nucleotide non-indel differences. The network  
500 analyses were performed using igraph implemented in R  
501 (<http://igraph.sourceforge.net/index.html>). The distribution of mismatches within a single  
502 network cluster were determined by aligning the sequence representing the largest vertex with  
503 the sequences to which it is connected and the positions of mismatches were determined  
504 along the sequences. Two-sided t-tests were performed in R.

505

#### 506 **Diversity measure calculations**

507 The Gini index was calculated by ordering the cluster sizes from largest to smallest  
508 and creating a cumulative frequency distribution, where  $R = \{r_1, r_2, \dots, r_n\}$ ,  $r_i$  is the  
509 cumulative size of the all the largest clusters until the  $i^{\text{th}}$  largest cluster and normalized such  
510 that  $r_n = 1$ . The Gini index is  $Gini\ index\ (g) = \sum_{i=1}^N \frac{(r_i - (i/N))}{N}$ , where  $N$  is the number of  
511 clusters (Morrow 1977).

512

#### 513 **Estimation of cluster sizes due to sequencing error**

514 The Poisson distribution can estimate the expected number of reads containing  $i$   
515 errors from the (central) vertex of size  $n$  reads, given an estimated error rate. The expected  
516 number of sequences with  $i$  errors is  $n.p_i$ , where  $p_i = P(X = i) = \frac{\lambda e^{-\lambda}}{i!}$ , and  $\lambda$  is the expected  
517 number of mutations per read. A cluster is defined as a set of interconnected vertices, in  
518 which edges are generated between vertices that differ by a single base. A vertex  $v$  is only

519 included in a cluster when the minimum distance from  $v$  to any of the sequences in the cluster  
 520 containing the central vertex is one. Thus, all the sequencing errors at  $i=l$  generate vertices  
 521 that have edges connecting to the central vertex. At  $i > l$ , a vertex with set of mutations  $M_x$   
 522 will be connected to the cluster only if there exists a vertex in the cluster with a set of  
 523 mutations  $M_y$  such that  $\left| \frac{M_x}{M_y} \right| = |\{x \in M_x | x \notin M_y\}| = 1$  (i.e. there is only one mutation in  
 524  $M_x$  that is not in  $M_y$ ). Therefore the probability of vertices due to  $i$  sequencing errors is  
 525 estimated by drawing  $S[n, i]$  samples from a multinomial distribution, for which the  
 526 probability of the possible vertices that could connect to the cluster is given by  $S[n, i] =$   
 527  $\prod_{j=1}^{i-1} \frac{E[n, j-1]}{l} \cdot p_i$ , where  $l$  is the length of the sequence and  $E[n, j]$  is the estimated number of  
 528 vertices that are in the cluster which are at distance of  $j$  from the central node. Here we draw  
 529 1000 independent samples from the multinomial distribution to estimate the average number  
 530 of vertices at distances  $i$  from the central vertex, and therefore the cluster size due to  
 531 sequencing error can be estimated by summing over the expected number of vertices at all  $i$ ,  
 532  $1 \leq i \leq \infty$ .

533

#### 534 **Temporal evolution of dominant clones:**

535 Sequences from the dominant clusters that were observed at least 6 times for the two  
 536 samples underwent a multiple alignment using the ClustalW2 algorithm  
 537 ([www.ebi.ac.uk/Tools/clustalw2/index.html](http://www.ebi.ac.uk/Tools/clustalw2/index.html)) with default parameters. A phylogenetic tree  
 538 was fitted using unrooted parsimony methods implemented in R ([http://cran.r-](http://cran.r-project.org/web/packages/phangorn/)  
 539 [project.org/web/packages/phangorn/](http://cran.r-project.org/web/packages/phangorn/)). Model tests were performed on different substitution  
 540 models, for which JC+G+I substitution model was found to be optimum and thus used here.  
 541 1000 bootstrap samples of individual nucleotides in the multiple alignment was used to assess  
 542 the reproducibility of the phylogenetic trees. The proportional difference in expression,  $diff(i)$ ,  
 543 of a given sequence  $i$  between month 0 and month 3 was calculated by difference in  
 544 expression between the two time points and divided by the sum of the expression of the  
 545 sequence over both times:

$$Diff(i) = \frac{E_i(t = 3) - E_i(t = 0)}{E_i(t = 3) + E_i(t = 0)}, \quad -1 \leq Diff(i) \leq 1$$

546 where  $E_i(t = 0)$  and  $E_i(t = 3)$  is the expression of sequence  $i$  at 0 months and 3 months

547 respectively.

548

549

550 **DATA ACCESS**

551 The IgH sequences discussed can be found as accession number ERP002120 in the European  
552 Nucleotide Archive (ENA).

553

554 **ACKNOWLEDGEMENTS**

555 This work was supported by the Wellcome Trust. We would like to thank the Cambridge Cancer Trials  
556 Centre and nurse specialists Gwyn Stafford, Rosie Tween, Lisa Walbridge and Joanna Baxter, and the  
557 patients and staff of Addenbrooke's Haematology Translational Research Laboratory and Cambridge  
558 Blood and Stem Cell Biobank which is supported by the BRC.

559

560 **AUTHORSHIP**

561 Contribution: R.J.M.B-R., A.L.P. and P.K. designed the study; R.J.M.B-R. performed experiments and  
562 analyzed the data; G.A.F. and G.S.V. provided patient samples; R.R. performed the 454 sequencing;  
563 B.J.H, G.A.F. and G.S.V. provided advice for the project; R.J.M.B-R., A.L.P. and P.K. wrote the paper,  
564 and all authors reviewed and approved the manuscript.

565 Conflict-of-interest disclosure: The authors declare no competing financial interests.

566

567

568 **References**

569

570 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* **215**(3): 403-410.571 Arber DA. 2000. Molecular diagnostic approach to non-Hodgkin's lymphoma. *J Mol Diagn* **2**(4): 178-190.572 Arnaout R, Lee W, Cahill P, Honan T, Sparrow T, Weiland M, Nusbaum C, Rajewsky K, Koralov SB. 2011. High-resolution description of antibody heavy-chain repertoires in humans. *PloS one* **6**(8): e22365.573 Bagnara D, Callea V, Stelitano C, Morabito F, Fabris S, Neri A, Zanardi S, Ghiotto F, Ciccone E, Grossi CE et al. 2006. IgV gene intraclonal diversification and clonal evolution in B-cell chronic lymphocytic leukaemia. *British journal of haematology* **133**(1): 50-58.574 Batrak V, Blagodatski A, Buerstedde JM. 2011. Understanding the immunoglobulin locus specificity of hypermutation. *Methods Mol Biol* **745**: 311-326.575 Ben-Hamo R, Efroni S. 2011. The whole-organism heavy chain B cell repertoire from Zebrafish self-organizes into distinct network features. *BMC Syst Biol* **5**: 27.576 Benichou J, Ben-Hamo R, Louzoun Y, Efroni S. 2012. Rep-Seq: uncovering the immunological repertoire through next-generation sequencing. *Immunology* **135**(3): 183-191.577 Bertoli D. 1997. Rapid amplification of cDNA ends. *Methods Mol Biol* **67**: 233-238.578 Boyd SD, Gaeta BA, Jackson KJ, Fire AZ, Marshall EL, Merker JD, Maniar JM, Zhang LN, Sahaf B, Jones CD et al. 2010. Individual variation in the germline Ig gene repertoire inferred from variable region gene rearrangements. *Journal of immunology* **184**(12): 6986-6992.579 Boyd SD, Marshall EL, Merker JD, Maniar JM, Zhang LN, Sahaf B, Jones CD, Simen BB, Hanczaruk B, Nguyen KD et al. 2009. Measurement and clinical monitoring of human lymphocyte clonality by massively parallel VDJ pyrosequencing. *Sci Transl Med* **1**(12): 12ra23.580 Brezinschek HP, Brezinschek RI, Lipsky PE. 1995. Analysis of the heavy chain repertoire of human peripheral B cells using single-cell polymerase chain reaction. *Journal of immunology* **155**(1): 190-202.581 Bruggemann M, White H, Gaulard P, Garcia-Sanz R, Gameiro P, Oeschger S, Jasani B, Ott M, Delsol G, Orfao A et al. 2007. Powerful strategy for polymerase chain reaction-based clonality assessment in T-cell malignancies Report of the BIOMED-2 Concerted Action BHM4 CT98-3936. *Leukemia* **21**(2): 215-221.582 Caligaris-Cappio F, Ghia P. 2008. Novel insights in chronic lymphocytic leukemia: are we getting closer to understanding the pathogenesis of the disease? *J Clin Oncol* **26**(27): 4497-4503.583 Campbell PJ, Pleasance ED, Stephens PJ, Dicks E, Rance R, Goodhead I, Follows GA, Green AR, Futreal PA, Stratton MR. 2008. Subclonal phylogenetic structures in cancer revealed by ultra-deep sequencing. *Proc Natl Acad Sci U S A* **105**(35): 13081-13086.584 Carulli G, Marini A, Ciancia EM, Bruno J, Vignati S, Lambelet P, Cannizzo E, Ottaviano V, Galimberti S, Caracciolo F et al. 2011. Discordant lymphoma consisting of splenic mantle cell lymphoma and marginal zone lymphoma involving the bone marrow and peripheral blood: a case report. *J Med Case Rep* **5**: 476.585 Dimitrov DS. 2010. Therapeutic antibodies, vaccines and antibodyomes. *mAbs* **2**(3): 347-356.586 Dorner T, Brezinschek HP, Foster SJ, Brezinschek RI, Farner NL, Lipsky PE. 1998. Delineation of selective influences shaping the mutated expressed human Ig heavy chain repertoire. *Journal of immunology* **160**(6): 2831-2841.587 Evans PAS, Pott C, Groenen PJTA, Salles G, Davi F, Berger F, Garcia JF, van Krieken JHJM, Pals S, Kluin P et al. 2007. Significantly improved PCR-based clonality testing in B-cell malignancies by use of multiple immunoglobulin gene targets. Report of the BIOMED-2 Concerted Action BHM4-CT98-3936. *Leukemia* **21**(2): 207-214.588 Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve LL, Gibbs RA, Belmont JW, Boudreau A, Hardenbol P, Leal SM et al. 2007. A second generation human haplotype map of over 3.1 million SNPs. *Nature* **449**(7164): 851-861.589 Freeman JD, Warren RL, Webb JR, Nelson BH, Holt RA. 2009. Profiling the T-cell receptor beta-chain repertoire by massively parallel sequencing. *Genome Res* **19**(10): 1817-1824.590 Ghia P, Caligaris-Cappio F. 2006. The origin of B-cell chronic lymphocytic leukemia. *Semin Oncol* **33**(2): 150-156.

591 Harris S, Bruggemann M, Groenen PJTA, Schuurin A, Langerak AW, Hodges E. 2012. Clonality analysis in lymphoproliferative disease using the

592

- 626 BIOMED-2 multiplex PCR protocols: experience from the EuroClonality  
627 group EQA scheme. *Journal of Hematopathology*: 91-98.
- 628 Hayes GM, Busch R, Voogt J, Siah IM, Gee TA, Hellerstein MK, Chiorazzi N, Rai KR, Murphy EJ.  
629 2010. Isolation of malignant B cells from patients with chronic lymphocytic leukemia (CLL)  
630 for analysis of cell proliferation: Validation of a simplified method suitable for multi-center  
631 clinical studies. *Leukemia Res* **34**(6): 809-815.
- 632 Hsi ED, Hoeltge G, Tubbs RR. 2000. Biclinal chronic lymphocytic leukemia. *Am J Clin Pathol*  
633 **113**(6): 798-804.
- 634 Jager U, Fridrik M, Zeitlinger M, Heintel D, Hopfinger G, Burgstaller S, Mannhalter C, Oberaigner W,  
635 Porpaczy E, Skrabs C et al. 2012. Rituximab serum concentrations during immuno-  
636 chemotherapy of follicular lymphoma correlate with patient gender, bone marrow infiltration  
637 and clinical response. *Haematologica* **97**(9): 1431-1438.
- 638 Jung D, Giallourakis C, Mostoslavsky R, Alt FW. 2006. Mechanism and control of V(D)J  
639 recombination at the immunoglobulin heavy chain locus. *Annu Rev Immunol* **24**: 541-570.
- 640 Kelly LM, Yu JC, Boulton CL, Apatira M, Li J, Sullivan CM, Williams I, Amaral SM, Curley DP,  
641 Duclos N et al. 2002. CT53518, a novel selective FLT3 antagonist for the treatment of acute  
642 myelogenous leukemia (AML). *Cancer Cell* **1**(5): 421-432.
- 643 Krause JC, Tsiabane T, Tumphey TM, Huffman CJ, Briney BS, Smith SA, Basler CF, Crowe JE, Jr.  
644 2011. Epitope-specific human influenza antibody repertoires diversify by B cell intraclonal  
645 sequence divergence and interclonal convergence. *Journal of immunology* **187**(7): 3704-3711.
- 646 Latchman D. 2005. Gene Regulation (Advanced Texts)
- 647 Lefranc MP, Giudicelli V, Ginestoux C, Jabado-Michaloud J, Folch G, Bellahcene F, Wu Y, Gemrot E,  
648 Brochet X, Lane J et al. 2009. IMGT, the international ImMunoGeneTics information system.  
649 *Nucleic Acids Res* **37**(Database issue): D1006-1012.
- 650 Lev A, Simon AJ, Bareket M, Bielorai B, Hutt D, Amariglio N, Rechavi G, Somech R. 2012. The  
651 kinetics of early T and B cell immune recovery after bone marrow transplantation in RAG-2-  
652 deficient SCID patients. *PloS one* **7**(1): e30494.
- 653 Lin MM, Zhu M, Scharff MD. 1997. Sequence dependent hypermutation of the immunoglobulin heavy  
654 chain in cultured B cells. *Proc Natl Acad Sci U S A* **94**(10): 5284-5289.
- 655 Logan AC, Gao H, Wang C, Sahaf B, Jones CD, Marshall EL, Buno I, Armstrong R, Fire AZ,  
656 Weinberg KI et al. 2011. High-throughput VDJ sequencing for quantification of minimal  
657 residual disease in chronic lymphocytic leukemia and immune reconstitution assessment.  
658 *Proceedings of the National Academy of Sciences of the United States of America* **108**(52):  
659 21194-21199.
- 660 Lukowsky A, Marchwat M, Sterry W, Gellrich S. 2006. Evaluation of B-cell clonality in archival skin  
661 biopsy samples of cutaneous B-cell lymphoma by immunoglobulin heavy chain gene  
662 polymerase chain reaction. *Leuk Lymphoma* **47**(3): 487-493.
- 663 Maletzki C, Jahnke A, Ostwald C, Klar E, Prall F, Linnebacher M. 2012. Ex-vivo clonally expanded B  
664 lymphocytes infiltrating colorectal carcinoma are of mature immunophenotype and produce  
665 functional IgG. *PloS one* **7**(2): e32639.
- 666 McClure RF, Kaur P, Pagel E, Ouillette PD, Holtegaard CE, Treptow CL, Kurtin PJ. 2006. Validation  
667 of immunoglobulin gene rearrangement detection by PCR using commercially available  
668 BIOMED-2 primers. *Leukemia* **20**(1): 176-179.
- 669 Morrow JS. 1977. Toward a more normative assessment of maldistribution: the Gini Index. *Inquiry*  
670 **14**(3): 278-292.
- 671 Rieben R, Frauenfelder A, Nydegger UE. 1996. Spectrotype analysis of human ABO antibodies:  
672 evidence for different clonal heterogeneity of IgM, IgG, and IgA antibody populations. *Vox*  
673 *Sang* **70**(2): 104-111.
- 674 Sanchez ML, Almeida J, Gonzalez D, Gonzalez M, Garcia-Marcos MA, Balanzategui A, Lopez-Berges  
675 MC, Nomdedeu J, Vallespi T, Barbon M et al. 2003. Incidence and clinicobiologic  
676 characteristics of leukemic B-cell chronic lymphoproliferative disorders with more than one  
677 B-cell clone. *Blood* **102**(8): 2994-3002.
- 678 Sandberg Y, van Gastel-Mol EJ, Verhaaf B, Lam KH, van Dongen JJ, Langerak AW. 2005. BIOMED-  
679 2 multiplex immunoglobulin/T-cell receptor polymerase chain reaction protocols can reliably  
680 replace Southern blot analysis in routine clonality diagnostics. *J Mol Diagn* **7**(4): 495-503.
- 681 Satoh M, Akizuki M, Yamagata H, Nakayama S, Homma M. 1996. Restricted heterogeneity and  
682 changing spectrotypes in autoantibodies to La/SS-B. *Autoimmunity* **24**(4): 229-236.
- 683 Sayala HA, Rawstron AC, Hillmen P. 2007. Minimal residual disease assessment in chronic  
684 lymphocytic leukaemia. *Best Pract Res Clin Haematol* **20**(3): 499-512.
- 685 Schatz DG, Swanson PC. 2010. V(D)J Recombination: Mechanisms of Initiation. *Annu Rev Genet*.

- 686 Schliep KP. 2011. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**(4): 592-593.
- 687 Sproul AM, Goodlad JR. 2012. Clonality testing of cutaneous lymphoid  
688 infiltrates: practicalities, pitfalls and potential uses. *Journal of*  
689 *Hematopathology*: 69-82.
- 690 Stamatopoulos K, Kosmas C, Stavroyianni N, Loukopoulos D. 1996. Evidence for immunoglobulin  
691 heavy chain variable region gene replacement in a patient with B cell chronic lymphocytic  
692 leukemia. *Leukemia* **10**(9): 1551-1556.
- 693 Stewart JJ, Lee CY, Ibrahim S, Watts P, Shlomchik M, Weigert M, Litwin S. 1997. A Shannon entropy  
694 analysis of immunoglobulin and T cell receptor. *Molecular immunology* **34**(15): 1067-1082.
- 695 Tonegawa S. 1983. Somatic generation of antibody diversity. *Nature* **302**(5909): 575-581.
- 696 van Dongen JJ, Langerak AW, Bruggemann M, Evans PA, Hummel M, Lavender FL, Delabesse E,  
697 Davi F, Schuurin E, Garcia-Sanz R et al. 2003. Design and standardization of PCR primers  
698 and protocols for detection of clonal immunoglobulin and T-cell receptor gene recombinations  
699 in suspect lymphoproliferations: report of the BIOMED-2 Concerted Action BMH4-CT98-  
700 3936. *Leukemia* **17**(12): 2257-2317.
- 701 van Krieken JH, Langerak AW, Macintyre EA, Kneba M, Hodges E, Sanz RG, Morgan GJ, Parreira A,  
702 Molina TJ, Cabecadas J et al. 2007. Improved reliability of lymphoma diagnostics via PCR-  
703 based clonality testing: report of the BIOMED-2 Concerted Action BHM4-CT98-3936.  
704 *Leukemia* **21**(2): 201-206.
- 705 Varadarajan N, Julg B, Yamanaka YJ, Chen H, Ogunniyi AO, McAndrew E, Porter LC, Piechocka-  
706 Trocha A, Hill BJ, Douek DC et al. 2011. A high-throughput single-cell analysis of human  
707 CD8(+) T cell functions reveals discordance for cytokine secretion and cytotoxicity. *J Clin Invest*  
708 **121**(11): 4322-4331.
- 709 Vargas RL, Felgar RE, Rothberg PG. 2008. Detection of clonality in lymphoproliferations using PCR  
710 of the antigen receptor genes: Does size matter? *Leukemia Res* **32**(2): 335-338.
- 711 Volkheimer AD, Weinberg JB, Beasley BE, Whitesides JF, Gockerman JP, Moore JO, Kelsoe G,  
712 Goodman BK, Levesque MC. 2007. Progressive immunoglobulin gene mutations in chronic  
713 lymphocytic leukemia: evidence for antigen-driven intraclonal diversification. *Blood* **109**(4):  
714 1559-1567.
- 715 Wang C, Mitsuya Y, Gharizadeh B, Ronaghi M, Shafer RW. 2007. Characterization of mutation  
716 spectra with ultra-deep pyrosequencing: application to HIV-1 drug resistance. *Genome Res*  
717 **17**(8): 1195-1201.
- 718 Warren RL, Freeman JD, Zeng T, Choe G, Munro S, Moore R, Webb JR, Holt RA. 2011. Exhaustive  
719 T-cell repertoire sequencing of human peripheral blood samples reveals signatures of antigen  
720 selection and a directly measured repertoire size of at least 1 million clonotypes. *Genome Res*  
721 **21**(5): 790-797.
- 722 Weinstein JA, Jiang N, White RA, 3rd, Fisher DS, Quake SR. 2009. High-throughput sequencing of  
723 the zebrafish antibody repertoire. *Science* **324**(5928): 807-810.
- 724 Williamson AR, Salaman MR, Kreth HW. 1973. Microheterogeneity and allomorphy of proteins.  
725 *Ann N Y Acad Sci* **209**: 210-224.
- 726 Woof JM, Burton DR. 2004. Human antibody-Fc receptor interactions illuminated by crystal  
727 structures. *Nature reviews Immunology* **4**(2): 89-99.
- 728 Wu YC, Kipling D, Leong HS, Martin V, Ademokun AA, Dunn-Walters DK. 2010. High-throughput  
729 immunoglobulin repertoire analysis distinguishes between human IgM memory and switched  
730 memory B-cell populations. *Blood* **116**(7): 1070-1078.
- 731  
732  
733

734 **FIGURE LEGENDS:**

735 **Figure 1. Sequencing of B-cell receptor repertoires. A)** Representation of the genomic  
736 rearrangement process during V-D-J recombination to generate the heavy chain B-cell  
737 receptor. B-cell receptor amplification was performed by reverse transcription on total RNA  
738 by single J region primer, and subsequent multiplex PCR amplification. **B)** The percentage of  
739 reads corresponding to the highest expressed B-cell receptor sequence for each sample,  
740 separated into sample type: healthy individuals, chronic lymphocytic leukemia patients (CLL)  
741 and human lymphoblastoid cell lines (LCL). Two-sided t-tests was performed between the  
742 sample subsets, with the p-values indicated above. **C)** Percentage of sequences shared  
743 between runs for technical repeats for a) the RT-PCR and re-sequencing (RT-PCR repeats,  
744 green bars) and b) the 454 sequencing from the same RT-PCR product (sequencing repeats,  
745 purple bars). For each sample, two repeats were performed and the percentage of reads shared  
746 between the repeats is shown (each repeat is compared to the other so two bars are shown per  
747 sample).

748

749 **Figure 2. B-cell receptor repertoires from different samples. A)** Schematic diagram  
750 showing the method by which the sequencing networks are generated: each vertex represents  
751 a unique sequence, where the relative size of the vertex is proportional to the number of 454  
752 sequencing reads that were identical to the vertex sequence. Edges are created between  
753 vertices that differ by one base (indel or substitution). The vertex colors correspond to the  
754 relative abundance of the corresponding sequences, where red, orange and yellow indicates  
755 observation of a sequence in >90%, between 40-90% and <40% of the reads in the sample  
756 respectively. **B)** Comparison of BCR sequence networks between **i)** a typical LCL sample and  
757 **ii)** a typical healthy individual. **C)** BCR sequence networks of CLL patients with **i)** extensive  
758 clonal enlargement and **ii)** limited clonal expansion. **D)** BCR sequence networks of CLL  
759 patient 5 showing expansion of two dominant clusters. **E)** Networks generated from  
760 sequencing dataset from *Boyd et al.* (Boyd et al. 2009) of **i)** healthy donor 1, **ii)** patient 2 with  
761 follicular lymphoma (FL) and **iii)** patient 3 with FL and small lymphocytic lymphoma (SLL).

762

763 **Figure 3. Measures differentiating between B-cell receptor populations. A)** Cluster Gini  
764 Index plotted against vertex Gini Index for thirteen healthy individual samples, eleven chronic  
765 lymphocytic leukemia (CLL), and eight human lymphoblastoid cell line (LCL) samples. Point  
766 (a) corresponds with healthy individual 10. The red box and grey dashed box distinguish  
767 between the regions occupied between diverse and clonal populations respectively. **B)** The  
768 second maximum cluster sizes plotted against the maximum cluster sizes. The red, grey  
769 dashed and black solid boxes distinguish between the regions occupied between unexpanded  
770 populations, monoclonal expanded populations and biclonally expanded populations  
771 respectively. **C)** B-cell receptor networks for the titration of a chronic lymphocytic leukemia  
772 clonal sample into healthy peripheral blood from the dataset from *Boyd et al.* (Boyd et al.  
773 2009) and **D)** the corresponding number of reads corresponding to the leukemic clone (green)  
774 and the maximum cluster size of each dilution (grey). The solid horizontal line shows the  
775 mean maximum cluster size for healthy individuals from this dataset (0.52% of total reads),  
776 and the dashed horizontal lines show the mean +/- standard deviation of maximum cluster  
777 size for healthy individuals for this dataset. **E)** Correlation between the Gini Index and the  
778 length of time since chronic lymphocytic leukemia (CLL) diagnosis for each patient in our  
779 dataset, with corresponding  $R^2$ -value.

780

781 **Figure 4. B-cell leukemic clonal evolution. A)** The B-cell sequence networks for patient A  
782 with chronic lymphocytic leukemia for samples i) prior to and ii) after second cycle of  
783 Chlorambucil treatment separated by 1 month with corresponding white blood cell counts. **B)**  
784 All sequences from the dominant clusters from both temporal samples were used to generate a  
785 composite network, and the differential frequencies at each time point is indicated by the  
786 relative vertex sizes. **C)** Correlation between the proportional frequencies of each unique  
787 BCR within the dominant clones of patient A with corresponding R-value and linear  
788 regression equation. **D)** An unrooted maximum parsimony tree was generated showing the  
789 relationships between sequences that were observed at least 6 times between the pre and post-

790 treatment samples, where the branch lengths are proportional to the number of varying bases  
791 (evolutionary distance). The same tree with the tip colors showing the relative difference in  
792 sequence abundance between the different time points, where green indicates observation of  
793 sequence primarily in the pre-treatment sample, blue indicates predominant observations in  
794 the post-treatment sample, and white indicates no change in frequency. Clones 1 and 2 refer  
795 to examples of BCRs observed only in the pre or post-treatment samples respectively.  
796

797 **Tables and figures:**

798

799 **Table 1. Sample information.**

<b>Sample</b>	<b>Patient type</b>	<b>Age, years</b>	<b>Gender</b>	<b>Time since CLL diagnosis, years</b>
CLL 1	CLL	77	Male	7
CLL 2	CLL	58	Male	2
CLL 3	CLL	78	Male	1.5
CLL 4	CLL+HCC	77	Male	2.5
CLL 5	CLL	59	Female	1.25
CLL 6	CLL	67	Male	2
CLL7	CLL	69	Male	13
CLL 8	CLL	64	Male	4.5
CLL 9	CLL	77	Male	5.25
CLL 10	CLL	81	Male	8
CLL 11	CLL	81	Male	10
Healthy 1	Age matched control 1	74	Female	-
Healthy 2	Age matched control 2	62	Female	-
Healthy 3	Age matched control 3	75	Female	-
Healthy 4	Age matched control 4	67	Female	-
Healthy 5	Age matched control 5	68	Female	-
Healthy 6	Healthy 6	55	Male	-
Healthy 7	Healthy 7	23	Male	-
Healthy 8	Healthy 8	23	Male	-
Healthy 9	Healthy 9	25	Male	-
Healthy 10	Healthy 10	24	Female	-
Healthy 11	Healthy 11	24	Female	-
Healthy 12	Healthy 12	24	Female	-
Healthy 13	Healthy 13	24	Female	-

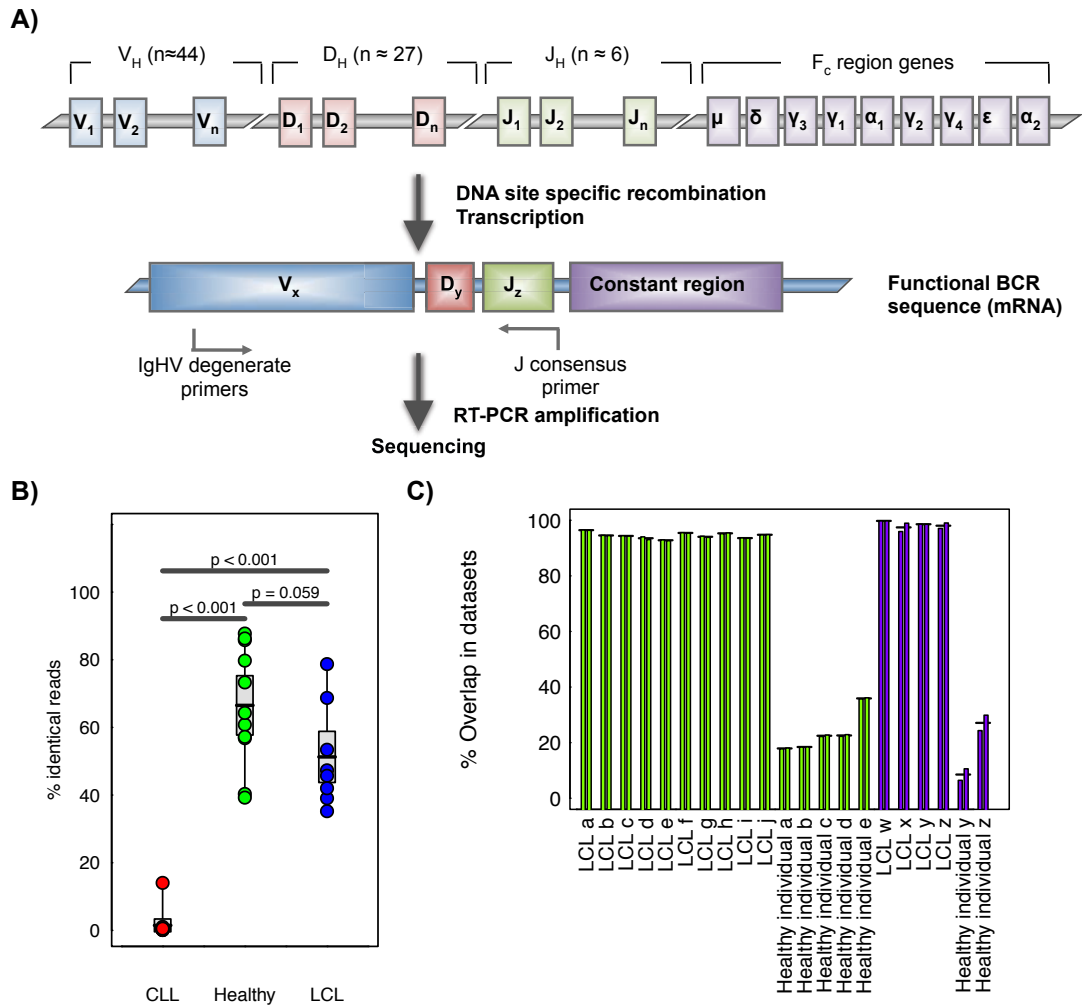
800 \* Abbreviations: LCL=Human lymphoblastoid cell line, CLL=chronic lymphocytic

801 leukemia, HCC=Hepatocellular carcinoma.

802

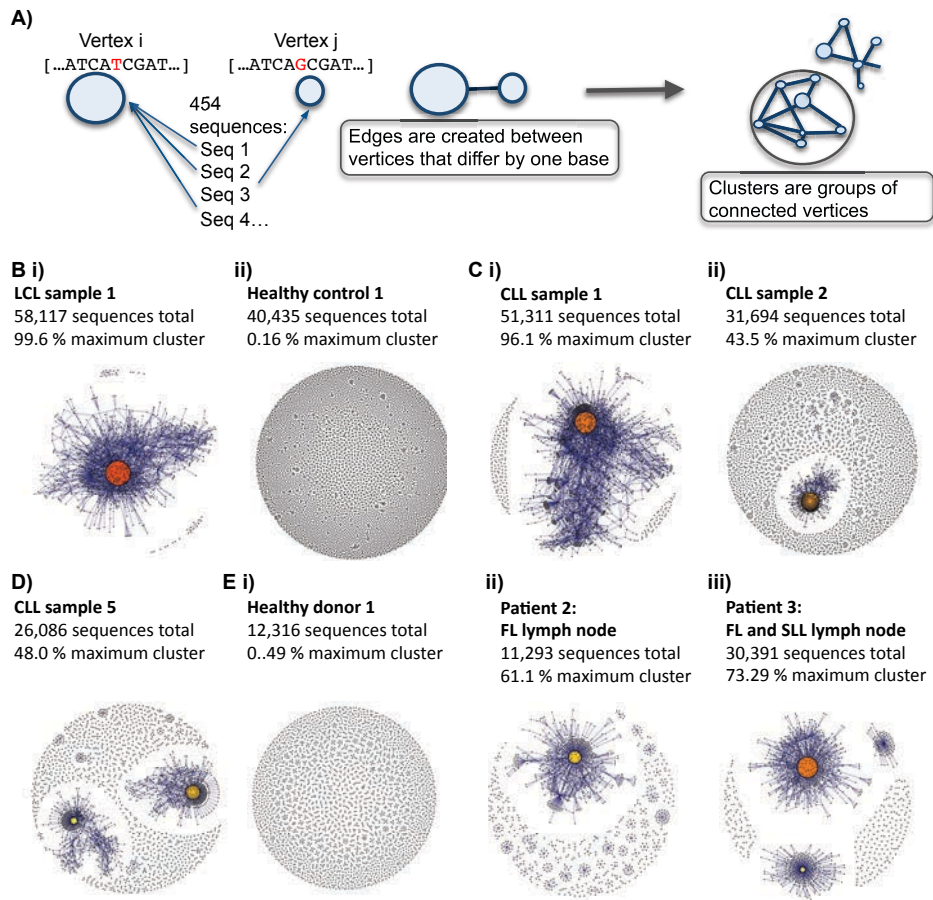
805  
806  
807

**Figure 1**



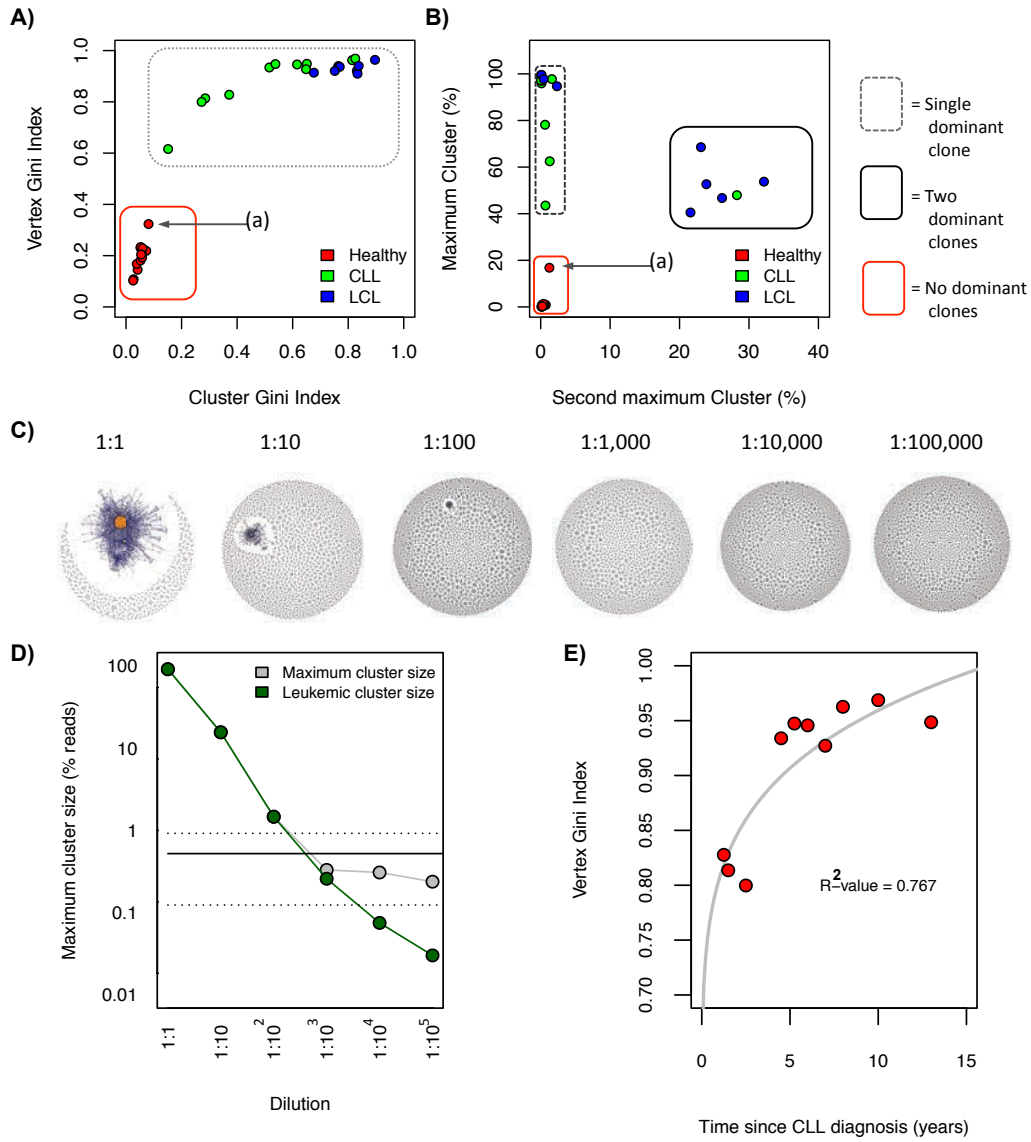
808  
809

810 **Figure 2**  
811



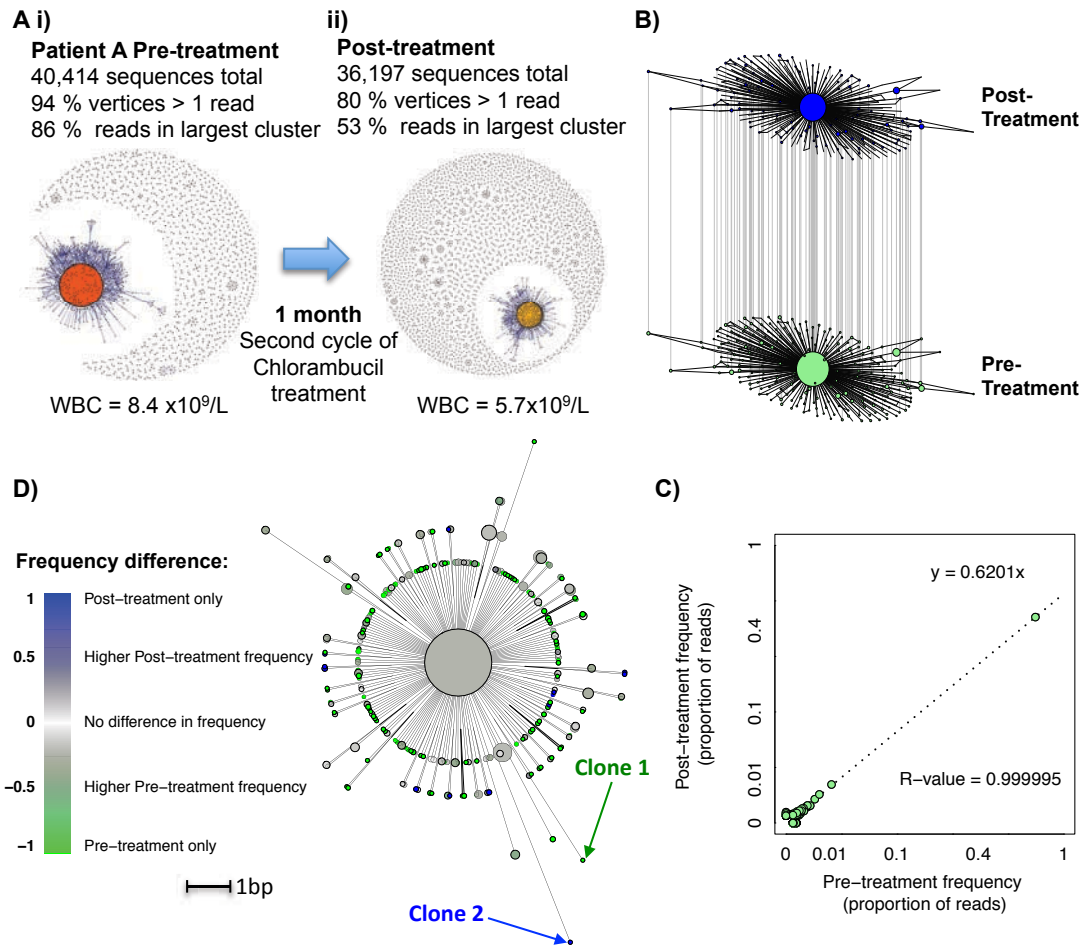
812  
813  
814

815 **Figure 3**  
816



817  
818

819 **Figure 4**  
 820  
 821



822  
 823