## 1 Supplementary methods:

2 **IGH leader sequencing**

3 We used 250 ng of buffy coat DNA as input for the leader-based multiplex PCR of IGHV-IGHD-IGHJ
4 rearrangements. After amplification, PCR products were purified on Agencourt AMPURE beads
5 (Beckman Coulter, Brea, CA, USA) and eluted using H2O. The purified product was checked for primer
6 dimers using the Agilent Tapestation (Agilent,Santa Clara, CA, USA). Paired-end sequencing was
7 performed using the MiSeq Reagent Kit v3 (2 × 300 bp) on the MiSeq Benchtop Sequencer (Illumina,
8 San Diego, CA, USA). PhiX was spiked-in at a 12% concentration to increase library diversity. The BcR
9 IGH repertoire was characterized using the ARResT/Interrogate immunoprofiler, an R/Shiny based tool
10 for in silico immunoprofiling developed by the Euroclonality-NGS working group. ARResT/Interrogate
11 annotates the variable (V), diversity (D) and joining (J) genes for each rearrangement and determines
12 the complementarity-determining region 3 (CDR3) by processing multiple IMGT/HighV-QUEST runs.
13 Clonotypes were computed as unique pairs of IGHV genes and CDR3 amino acid sequences within a
14 given sample. Only reads annotated as productive IGH rearrangements were included to prevent
15 underestimation of the dominant clonotype frequency of clones with bi-allelic rearrangements.
16 ARResT/Interrogate uses IMGT's definition of an unproductive rearrangements as an out-of-frame
17 junction, with one or more of the following: stop codon(s), frameshift mutation(s), defects in the
18 splicing sites and/or the regulatory element(s), unusual features such as translocation or gene fusion,
19 and/or changes of conserved amino acids demonstrated to lead to incorrect folding. QC metrics for
20 the sequencing data are shown in table S6. Out of the 277 samples (118 controls, 124 patients and 35
21 repeated samples) sequenced in this study, 8 samples did not pass the initial QC filter. Primary cause
22 of QC fail was the presence of primer dimers in the sample, introducing reads shorter than 50 nt. After
23 cleanup of primer dimers, a mean of 252,935 high-quality sequences were retrieved per sample.

24 **Ethical approval**

25 The study was approved by the local institutional medical ethical committee at the Erasmus MC
26 (protocol number MEC 2019-0484). The EPIC steering committee approved the use of the material for
27 the purpose of this study. All patients gave their written consent and the use of the material and data
28 in this study were approved by the IARC Ethics Committee. The study was performed in compliance
29 with the declaration of Helsinki.

30 **Case/Control matching**

31 Controls were matched on age, sex, EPIC center (and thus country) and blood draw date. Controls
32 were alive and without a cancer diagnosis (other than nonmelanoma skin cancer) at the time of the
33 diagnosis of CLL. For descriptive data see table S1.

34 **Stereotyped CLL subsets**

35 Stereotyped subsets were initially annotated through the ARResT/AssignSubsets tool and validated
36 through algorithms developed at INAB|CERTH (Thessaloniki, Greece). Stereotyped subsets were
37 defined by the following parameters: (1) usage of IGHV genes from the same phylogenetic clan, (2) a
38 minimum of 50% amino acid identity and 70% similarity within the heavy chain CDR3, (3) identical
39 heavy chain CDR3 length and, (4) identical offset of the shared amino acid pattern.[14] In the initial study

40 by Agathangelidis *et al*., a subset was defined as major if it represented at least 0.2% of their study
41 cohort, which amounted to 60 cases, subsets below this cutoff are referred to as minor subsets. IGL
42 light chain Sanger sequencing was performed using an IGLV3-21 specific forward primer and intron-
43 based reverse primers.

44 **Clinical data**

45 The label for CLL in the EPIC database was shared with SLL and distinction between these entities was
46 not possible based on the available information. Additional clinical data for 32/124 patients was
47 obtained by direct collaboration with one of the EPIC centers, (Umeå University, Sweden)
48 (**Supplementary figure 1,Table S2**). Sharing of the data and material was approved by the Swedish
49 Ethical Review Authority. All material and data received were anonymized and solely accessible for
50 researchers directly involved with the project. In total, 3/32 patients were diagnosed with SLL (9.4%)
51 and 26/32 patients were diagnosed with CLL (81%). The three remaining patients with a CLL/SLL label
52 in the EPIC database were reclassified as HC-MBL upon review of the clinical data. In general, it is
53 important to consider that 105/124 patients described in the current study were diagnosed with CLL
54 before 2008, when the diagnostic criterium for CLL was set at a persisting monoclonal B-cell count of
55 $5 \times 10^9$ cells/L. Before 2008, CLL Rai stage 0 was diagnosed based on an excess of $5 *10^9$ lymphocytes/L,
56 which means some cases we would now diagnose as HC-MBL would be included. As we do not have
57 detailed clinical data for all patients, retrospective reclassification of these patients was impossible.

58 **Sampling instances (Initial, longitudinal and diagnostic)**

59 Although the 242 samples were taken up over a period of 22 years before CLL diagnosis, blood
60 sampling over this period was not evenly distributed (**Table S3**). As a result, our findings become
61 increasingly uncertain as time to diagnosis increases. Realistically, no conclusions can be drawn earlier
62 than 18 years before diagnosis. For 22 patients repeated samples were available (N=35), as they also
63 contributed to the Northern Sweden Health and Disease Study. Samples were obtained by direct
64 collaboration with one of the EPIC centers, (Umeå University, Sweden). For 16 of these patients, a
65 diagnostic sample was available (**Supplementary figure 1BC**). For 15 out of the 22 patients, an
66 additional longitudinal pre-diagnostic sample was available. Of particular interest were the nine
67 patients who had both a diagnostic and multiple pre-diagnostic samples available.

68 **Statistics**

69 Frequency of a clonotype indicates the % productive reads of the total BcR IGH gene repertoire in the
70 patient sample. Comparison of dominant clonotype frequency between patients and controls was
71 done through two-tailed Wilcoxon–Mann–Whitney two-sample rank-sum test. For comparisons
72 between patients and controls, only the earliest sample available was used. Correlation in dominant
73 clonotype frequency vs. time to diagnosis in CLL patients was done using Spearman's rank correlation
74 and a line was fitted using local polynomial regression, also known as locally estimated scatterplot
75 smoothing (loess). All statistics and plotting was conducted in R (R Core Team, 2021). Kaplan-Meier
76 survival curves were plotted through the Survminer R package (Kassambra, 2017). Survival
77 distributions were compared by log-rank test. The linear mixed effects model for repeated samples
78 was fitted using the nlme R package (Pinheiro J, Bates D, DebRoy S, Sarkar D, R Core Team, 2021).
79 Linear mixed effects model included CLL clonotype frequency as outcome variable, with time to
80 diagnosis, age at blood draw, sex, mutational status, Rai stage at diagnosis and Binet stage as fixed

81 effects and patient ID as random effect. The model was fitted using maximum likelihood estimation
82 (MLE). The effect of time to diagnosis was allowed to differ between patients. Linear mixed effects
83 model goodness of fit was assessed by Likelihood Ratio Test. Addition of splines or interaction terms
84 did not improve the model. Significance of the fixed effects was determined by Wald test, with only
85 time to diagnosis having a significant effect on the CLL clonotype frequency ($P < 0.0001$), though this
86 may be due to the limited amount of patients with repeated samples. All statistical tests were two-
87 tailed.

88

89 Supplemental figure 1. **Overview of patient material and data.** A) Flowchart of the analysis, including patient
90 and control counts for each step. B) Chord diagram indicating the additional information available for the cohort.
91 C) Specification of additional clinical data and repeated samples. A total of 32 patients had clinical data available,
92 of whom 6 only had an additional longitudinal sample, 7 solely had an additional sample available at diagnosis
93 and 9 had both a longitudinal sample available and a matching diagnostic sample. For the remaining 10 patients
94 with clinical data no additional samples were available.

95 Supplemental figure 2. **Overall survival after CLL diagnosis.** A) Overall survival in years since CLL diagnosis for
96 CLL patients compared to controls. For controls, the date of CLL diagnosis of the matched CLL patient is used
97 instead. B) Overall survival in years after CLL diagnosis for CLL patients with a pre-diagnostic IGHV mutated
98 clonotype >2% of the total IGH gene repertoire compared to CLL patients with a pre-diagnostic IGHV unmutated
99 clonotype >2%. C) Overall survival in years after CLL diagnosis for CLL patients with a pre-diagnostic clonotype
100 >2% of the total IGH gene repertoire compared to CLL patients without a clonotype >2% of the IGH gene
101 repertoire.

102 Supplementary Table 1. **Descriptive table of the EPIC participants.**

103 Supplementary Table 2. **Clinical data of patients.**

104 Supplementary Table 3. **Dominant clonotype frequency over time for patients and controls.**
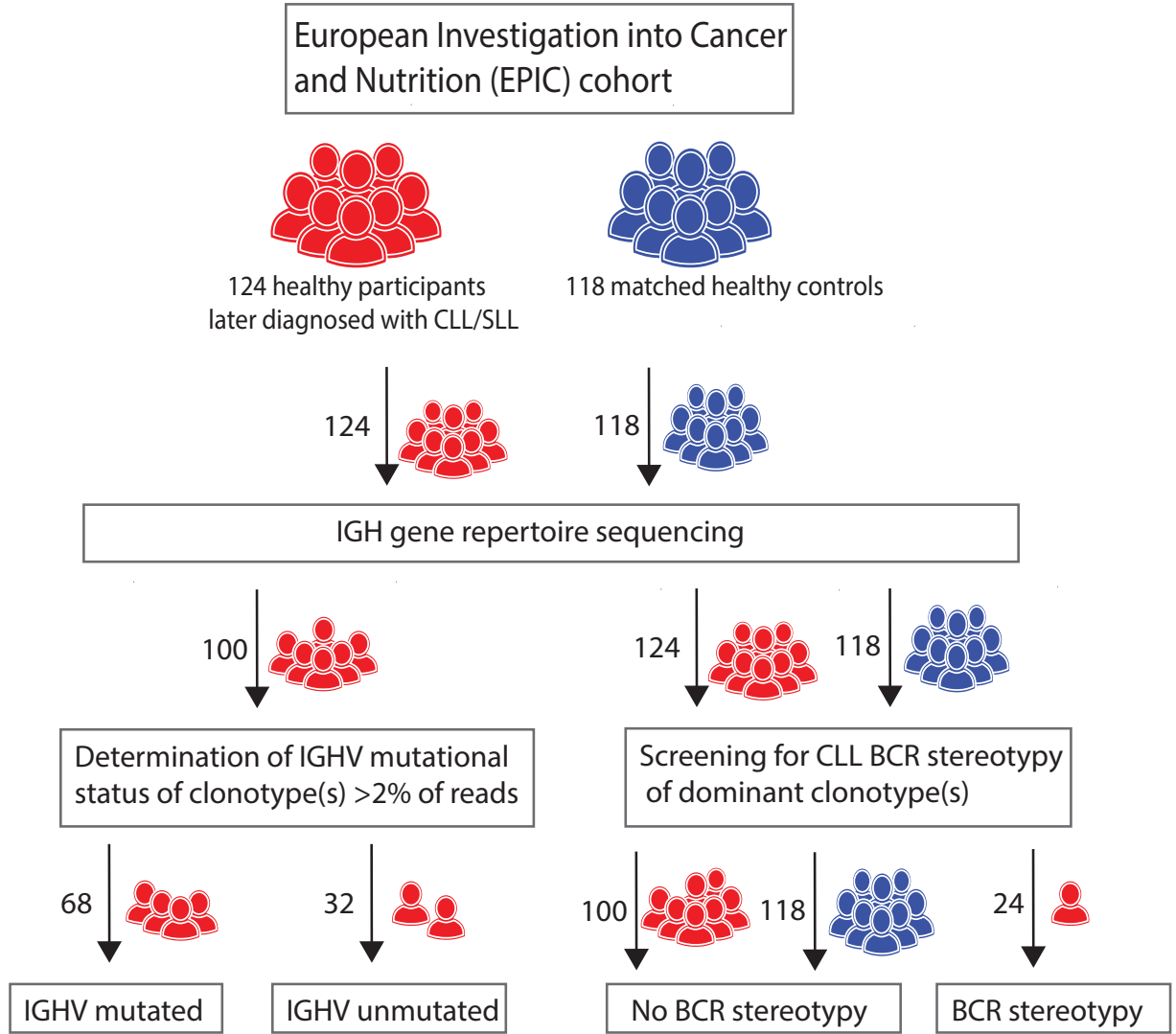
105 Supplementary Table 4. **Detailed overview of all CLL stereotyped subsets encountered in this study.**

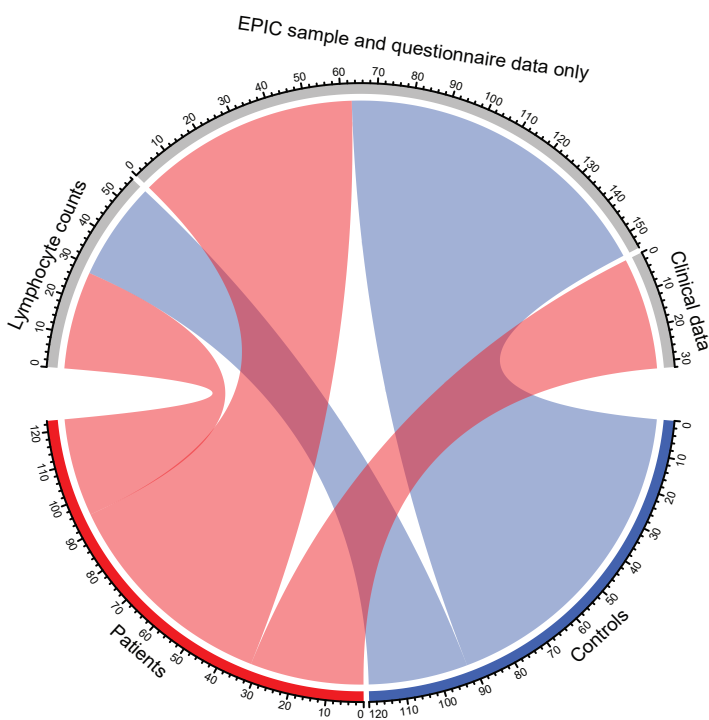106 Supplementary table 5. **Overview of patients with diagnostic material.**

107 Supplementary table 6. **NGS QC metrics.**

108

A

European Investigation into Cancer and Nutrition (EPIC) cohort

124 healthy participants later diagnosed with CLL/SLL

118 matched healthy controls

124 → 118 →

IGH gene repertoire sequencing

100 →

Determination of IGHV mutational status of clonotype(s) >2% of reads

124 → 118 →

Screening for CLL BCR stereotypy of dominant clonotype(s)

68 → IGHV mutated

32 → IGHV unmutated

100 → 118 → No BCR stereotypy

24 → BCR stereotypy

B

EPIC sample and questionnaire data only

Lymphocyte counts

Clinical data

Patients

Controls

C

Nr. of patients

Data and material available
Clinical data only
Diagnostic + Longitudinal
Diagnostic sample
Longitudinal sample(s)

Cases with clinical data

Supplemental figure 2