

1 **IgSeqR: a protocol for the identification, assembly, and characterization of full-**
2 **length tumor Immunoglobulin transcripts from unselected RNA sequencing**
3 **data**

4
5 Dean Bryant*,¹ Benjamin Sale*,¹ Giorgia Chiodin,¹ Dylan Tatterton,¹ Benjamin
6 Stevens¹, Alyssa Adlaon¹, Erin Snook¹, James Batchelor,^{1,2} Alberto Orfao,³ Francesco
7 Forconi^{1,4}

8 ¹Cancer Sciences, Faculty of Medicine, University of Southampton, Southampton, UK.

9 ²Clinical Informatics Research Unit, University of Southampton, Southampton, UK.

10 ³Cancer Research Center (IBMCC, USAL-CSIC), Cytometry Service (NUCLEUS),
11 Department of Medicine, Biomedical Research Institute of Salamanca (IBSAL),
12 University of Salamanca, Salamanca, Spain. ⁴Haematology Department, Cancer Care
13 Directorate, University Hospital Southampton NHS Trust, Southampton, UK.

14 *D.B and B.J.S have equally contributed

15

16 **Correspondence:** Francesco Forconi, Cancer B cell group, Cancer Sciences,
17 Somers Building, MP824, Tremona Road, Southampton, SO16 6YD, UK. **Email:**
18 f.forconi@soton.ac.uk. **Tel:** +44 (0)23 81205780

19

20

21

22

23

24

25

26

27

28

1 **Abstract**

2 Immunoglobulin (IG) gene analysis provides fundamental insight into B-cell receptor structure
3 and function. In B-cell tumors, it can inform the cell of origin and clinical outcomes. Its clinical
4 value has been established in the two types of chronic lymphocytic leukemia with unmutated
5 or mutated *IGHV* genes and is emerging in other B-cell tumors. The traditional PCR-based
6 techniques, which are labor-intensive, rely on the attainment of either a dominant sequence
7 or a small number of subclonal sequences and do not allow automated matching with the
8 clonal phenotypic features. Extraction of the expressed tumor IG transcripts using high-
9 throughput RNA sequencing (RNA-seq) can be faster and allow the collection of multiple
10 sequences matched with the transcriptome profile. Analytical tools are regularly sought to
11 increase the accuracy, depth, and speed of acquisition of the full *IGV-(IGD)-IGJ-IGC*
12 sequences and combine the *IG* characteristics with other RNA-seq data. We provide here a
13 user-friendly protocol for the rapid extraction, identification, and accurate determination of the
14 full (leader to constant region) tumor *IG* templated and non-templated transcript sequence
15 from RNA-seq. The derived amino acid sequences can be interrogated for their physico-
16 chemical characteristics and, in certain lymphomas, predict tumor glycan types occupying
17 acquired N-glycosylation sites. These features will then be available for association studies
18 with the tumor transcriptome. The resulting information can also help refine diagnosis,
19 prognosis, and potential therapeutic targeting in the most common lymphomas.

20 Word counts: 227 (max 250)

21

1 Introduction

2 The B-cell receptor (BCR) immunoglobulin (IG) glycoprotein is the defining functional
3 feature of a mature B cell, and *IG* gene analysis can provide fundamental insight into
4 the origin and behavior of a B-cell tumor [1, 2]. It is a Y-shaped dimer of 2 identical
5 heavy and light chains, with 2 main functional components. A variable region that
6 confers diversity to recognize different antigens and is unique to each B cell, and a
7 constant region with an effector function. *IG* diversity results from a series of genetic
8 recombinations at the *IG* heavy (*IGH*) and kappa (*IGK*) or lambda (*IGL*) light chain loci
9 during B cell development in the bone marrow before a naïve B cell exits to the
10 periphery (**Figure 1**). For the heavy chain, the recombinations are accompanied by
11 non-templated nucleotide additions/deletions at the junctions of one of ~51 *IGHV*, ~21
12 *IGHD*, and ~6 *IGHJ* genes in the complementarity-determining region 3 (*CDR3*)
13 forming the “fingerprint” of an individual B cell. Further variability is conferred by the
14 recombination of a *V* gene with a *J* gene at the *IGK* or *IGL* loci. Following antigen
15 encounter, naïve B-cells undergo class-switch recombination and somatic
16 hypermutation, typically in the presence of activation-induced cytidine deaminase
17 (AID), T cells, and cytokines, for affinity maturation in the germinal center (GC) and
18 differentiation in memory B cells or plasma cells [3]. The GC reaction involves
19 proliferation, which makes the B cells vulnerable to damage and transformation into
20 tumors. Tumor B cells preserve the *IG* sequence of the cell of origin. Therefore,
21 analysis of the *IG* sequences allows the identification of the stage of differentiation
22 reached by a B-cell before tumor transformation [4-6].

23 In chronic lymphocytic leukemia (CLL), *IG* analysis reveals two major types defined by
24 *IGHV* mutational status [5]. The CLL type with unmutated *IGHV* (U-CLL) derives from
25 pre-germinal center CD5⁺ B cells, while the CLL type with mutated *IGHV* (M-CLL)
26 appears to arise from post-follicular CD5⁺ B cells [7, 8]. Since the discovery that U-
27 CLL has a worse prognosis than M-CLL [9, 10], subsequent studies have
28 demonstrated that each type has a distinctive cellular origin, biology, (epi)genetics,
29 clinical prognosis, and response to therapy [5, 11]. *IGHV* gene analysis has become
30 an essential part of the diagnostic workup for any patient with CLL.

31 *IG* analysis also informs key tumor-specific features in certain lymphomas. In classic
32 follicular lymphoma (FL), the tumor *IG* acquires N-glycosylation sites (AGS), defined
33 by the asparagine-X-serine/threonine motif (where X is any amino acid except proline)
34 [12]. AGS in FL are typically in the CDRs of the variable region by somatic
35 hypermutation [13] and are occupied by tumor-specific oligomannose-type glycans
36 [14-16]. These atypical glycans are unique to the tumor B cell, are present on the entire
37 FL clone, and persist during the entire clonal history of FL through transformation into
38 diffuse large B-cell lymphoma (DLBCL), despite ongoing somatic hypermutation [14,
39 17].

40 The current gold standard for *IG* gene analysis is by Sanger sequencing. This
41 approach offers a highly accurate *IG* sequence but is time-consuming, labor-intensive,
42 and requires a dedicated experimental and analytical workflow on samples with
43 documented high tumor infiltration [18]. The increasing adoption of high-throughput
44 whole transcriptome RNA sequencing (RNA-seq) methods allows many tests to be
45 streamlined into a single experimental workflow. Through the application of
46 appropriate analytical pipelines, a single RNA-seq experiment can yield
47 comprehensive information on gene expression, isoform expression, single nucleotide
48 polymorphisms (SNPs), and larger structural variants [19].

1 RNA-seq can therefore be a better alternative to Sanger sequencing in *IG* gene
2 analysis. However, the intrinsic high variability of the non-templated CDR3 sequences
3 has been a challenge to the identification of the full *IG* sequence with the current RNA-
4 seq analytical workflows, which have involved mapping reads to a reference
5 transcriptome.

6 Here we describe IgSeqR (pronounced I-G-Seeker), a protocol for the reference-free
7 extraction, identification, and accurate determination of the full tumor *IG* transcript
8 sequence from unselected whole transcriptome RNA-seq data.

9 *Development of the protocol*

10 We first used IgSeqR to identify the tumor *IG* full transcripts in a cohort of 489 DLBCL
11 with RNA-seq data publically available [14]. The data were deposited in the National
12 Cancer Institute (NCI) Genomic Data Commons (accession phs001444.v1.p1) [20,
13 21]. The full *IGHV-IGHD-IGHJ* sequence rearrangements were identified from leader
14 to constant region with high confidence in 339 (69%) samples, from which we could
15 determine *IGHV*, *IGHD*, *IGHJ*, and *IGHC* use, homology to germline, and AGS
16 presence and location. Since we were interested in those cases with N-glycosylation
17 sites acquired by somatic hypermutation and no information was available on the
18 tumor purity of these samples, we investigated only the 307 samples with mutated
19 (<98% homology to germline) *IGHV* [14]. We found that the AGS were preferentially
20 in the EZB genetic subtype of the GC-B-cell-like (GCB) DLBCL. The majority of these
21 AGS were located in the CDR, in a fashion similar to FL. Following the generation of
22 F(ab) from the tumor-derived *IG* heavy and light chain sequences we documented that
23 the glycan structure occupying the AGS was location-dependent and that the
24 oligomannose-type glycans occupied the CDR-located sites only. We performed
25 correlation studies with the transcriptome profile and defined genes and gene sets
26 differentially expressed in samples with and without AGS. We performed correlations
27 with the clinical characteristics of the DLBCL. Interestingly, we found that AGS in the
28 EZB subtype conferred a poor prognosis, indicating that this approach for *IG* gene
29 analysis could be adopted to predict both glycan structure and response to
30 conventional therapies [14]. In the present study, we report the IgSeqR script while
31 validating its accuracy in primary CLL samples with matched *IG* heavy chain Sanger
32 and bulk RNA-seq data (deposited in ArrayExpress, accession E-MTAB-12017) [22].
33 IgSeqR is fully concordant with Sanger sequencing for *IGHV*, *IGHD*, and *IGHJ* allele
34 use and nucleotide sequence.

35 *Applications of the method*

36 IgSeqR is ideal for studies requiring high-quality base calls across the full sequence,
37 including the non-templated CDR3 region, of the *IG* heavy and light chains of any
38 mature B cell tumor. The protocol reduces the computational burden of *de novo*
39 assembly by pre-filtering redundant data and allows the identification of the dominant
40 nucleotide sequence of the *IG* heavy and light chains from leader to constant region
41 from RNA-seq data. Through the alignment to the most updated *IG* sequence
42 repertoires, currently IMGT/V-QUEST reference directory 202349-3, program version
43 3.6.2 at <http://www.imgt.org>, it is possible to obtain insights into *IGHV*, *IGHD*, *IGHJ*
44 heavy chain alleles, *IGKV*, *IGKJ* or *IGLV*, *IGLJ* light chain alleles, constant region class
45 and subclass, homology to germline, CDR1-3 and FR1-4 characteristics.

1 IgSeqR can also be applied to autoimmune and infectious diseases to identify
2 common patterns recurring in the polyclonal expansions (e.g. dominance and
3 characteristics of *IGHV1-69* in rheumatoid arthritis or influenza) [23, 24].

4 IgSeqR can also be used to generate F(ab)s [14] or improve strategies for vaccine
5 and antibody therapy development [25-27].

6 *Comparisons with other methods*

7 Compared to Sanger sequencing and existing RNAseq-based protocols, IgSeqR
8 increases the length of the transcript containing the full *IGV(-IGD)-IGJ* rearrangements
9 from leader to *IG* constant region (up to 2000 nucleotides).

10 It maintains the same level of accuracy as Sanger sequencing while improving the
11 chance of detecting a clonal sequence compared to a PCR-based approach,
12 particularly in lymphoma samples. *IG* sequencing of lymphoma samples by Sanger is
13 notoriously difficult and demands significant amounts of equipment and time,
14 particularly if subcloning approaches are necessary, to identify small cohorts of
15 patients [28-30]. In a cohort of 37 lymphomas with more than 10% tumor B cells in the
16 test sample by flow cytometry, PCR/direct Sanger sequencing successfully identified
17 a dominant *IG* rearrangement in only 11 (30%). By Cibersort estimation [31], 439
18 DLBCL samples from the NCI cohort had > 10% (tumor) B cells. IgSeqR identified the
19 tumor *IG* rearrangement in 319 (73%), a significantly superior frequency than Sanger
20 ($p < 0.0001$). However, IgSeqR was also successful in identifying the full *IG* sequence
21 in 20 of the 50 (40%) samples with <10% B cell purity, although the success rate was
22 lower compared to >10% ($p < 0.005$) (**Figure 2** and **Table S1**).

23 The experimental and analytical time to identify the sequences by Sanger was in
24 weeks, while it was in days for the IgSeqR approach. This suggests that IgSeqR is
25 dramatically efficient, offering a higher success rate in a shorter experimental and
26 analytical time compared to standard PCR and Sanger sequencing.

27 Although IgSeqR is currently not configured to build the *IGHC* sequence with contigs
28 spanning from CDR3 to the 3' end of the constant region allele used, the derived
29 transcripts recovered are generally sufficient to determine the *IGHC* class and
30 subclass with high confidence. This is another advantage compared to Sanger, where
31 individual isotypes can only be identified using isotype-specific primers.

32 Several tools have been developed for *IG* analysis from bulk and single-cell RNA-seq
33 [32-40] (**Table 1**), many of which preferentially rely on aligning RNA-seq reads to *IG*
34 reference sequences [33, 35-37]. MiXCR is widely adopted for immune profiling in
35 both academic and industrial settings [33]. It primarily uses the N-regions at the *IGV-*
36 (*IGD*)-*IGJ* junctions as a reference and identifies and quantifies the *IG* repertoire by
37 CDR3 diversity. However, it is less focused on the full length, and highly mutated *IGV-*
38 (*IGD*)-*IGJ* sequences may not be fully reconstructed. TRUST4 and IG_ID tools utilize
39 *de novo* transcriptome assembly. However, TRUST4 was initially designed for TCR,
40 rather than BCR, repertoire analysis [39]. The IG_ID tool can accurately produce full-
41 length BCR transcripts comparable to Sanger sequencing, but has an extended
42 processing time and generates large temporary files due to the *de novo* assembly of
43 the whole transcriptome, limiting its use for large-scale analyses [32].

44 We performed a direct comparison of IgSeqR with the MiXCR (v 4.3.2) or TRUST4
45 (v1.0.12) with the 18 CLL samples (**Tables 2** and **S2**).

1 MiXCR generated *IGH* transcripts for all 18 of the samples, but only 17 (94%) of these
2 spanned the full *IGHV-IGHD-IGHJ* rearrangement, and only 14 (78%) had 100%
3 identity with Sanger.

4 TRUST4 generated *IGH* transcripts from 17 (94%) of the samples, all of which were
5 fully concordant with Sanger. However, TRUST4 failed to identify the only case that
6 had a deletion of codon 66 of the *IGHV4-34* tumor sequence, possibly revealing a
7 limitation of TRUST4 in identifying insertions or deletions.

8 IgSeqR also produced the longest tumor transcripts, averaging a length of 2036
9 nucleotides, compared to 589 and 769 nucleotides by MiXCR and TRUST4
10 respectively. Notably, the majority (78%) of the IgSeqR transcripts were long enough
11 to cover the full *IGH* region from leader to the membrane domain of the constant region
12 with confidence, a feature not possible in the shorter transcripts generated by MiXCR
13 or TRUST4 (**Figure 3**).

14 When efficiency was assessed, IgSeqR took on average 1.18 seconds per nucleotide
15 assembled (s/nt) to complete, compared to 8.10 s/nt and 1.44 s/nt minutes by MiXCR
16 and TRUST4, respectively (**Table S3**).

17 Overall, IgSeqR obtained longer transcripts, was more efficient per nucleotide
18 assembled, and was more accurate than MiXCR and TRUST4.

19 *Experimental Design*

20 The experimental design of IgSeqR is divided into four key stages (**Figure 3**): (a) data
21 preprocessing, (b) *de novo* transcriptome assembly, (c) *IG* transcript selection and
22 quantification, and (d) *IG* transcript annotation and interpretation.

23 **Data preprocessing**

24 IgSeqR can use RNA-seq data in either BAM or FASTQ format. We have assessed
25 the quality of the RNA-seq data using FastQC [41], but alternative methods more
26 familiar to the operator can be used. Alignment of the data to a reference transcriptome
27 is performed using HISAT2 alignment tool, which employs a hierarchical indexing
28 strategy based on Burrows-Wheeler Transform [42]. If the input file has been
29 previously aligned, FASTQ reads must first be extracted from the alignment file (Step
30 1) before being supplied to HISAT2 (Step 2). It is problematic to map *IG* variable
31 genes, especially D and J to a reference transcriptome using short read RNA-seq data,
32 which results in many *IG*-derived reads being unmapped following alignment [32].
33 Therefore, following alignment, the resultant BAM file is filtered to extract reads which
34 align to specific *IG* associated genomic loci in addition to any unmapped reads.

35 **De Novo Assembly**

36 The Trinity software [43] is used for reference-free transcript reconstruction of the
37 reads associated with *IG* sequences. Trinity follows a three-step process: Inchworm,
38 Chrysalis, and Butterfly [43]. Inchworm builds initial contigs by assembling overlapping
39 k-mers from the short reads. Chrysalis constructs a De Bruijn graph using the
40 Inchworm contigs to represent connections between overlapping sequences and
41 identifies alternative splicing events. Butterfly decomposes the De Bruijn graph into
42 individual components representing distinct transcripts from the same gene. These
43 components are refined and merged to generate complete transcript sequences. The
44 filtered FASTQ files generated from the HISAT2 output are supplied to Trinity for *de*

1 *novo* transcript assembly, resulting in a FASTA file containing the assembled
2 transcripts (Step 4).

3 **IG Transcript Selection and Quantification**

4 To remove any non-IG associated transcripts assembled by Trinity, the transcripts in
5 the output FASTA file are aligned to reference IG databases using BLAST [44].
6 Reference FASTA IG sequences are concatenated to generate the databases (Step
7 5). Transcripts that align with an IG reference sequence are retained (Step 6) and
8 quantified using Kallisto [45], a tool that quantifies transcript abundance from RNA-
9 Seq data using pseudo-alignment instead of read alignment. A k-mer-based index is
10 built (Step 7) for quantification of the filtered transcripts using the FASTQ reads (Step
11 8).

12 **IG Transcript Annotation and Interpretation**

13 The most abundant transcripts are identified using the transcript quantification outputs
14 (Step 9) and passed through the IMGT/V-QUEST sequence alignment web tool [46],
15 benefiting from a comprehensive database of known germline IG alleles and
16 polymorphisms for functional annotations (Step 10). V-QUEST identifies and
17 annotates *IGHV-IGHD-IGHJ* and *IGKV-IGKJ* or *IGLV-IGLJ* rearrangements, detects
18 nucleotide mutations and insertions/deletions, and functionality. The annotated
19 transcripts are then manually reviewed to identify the tumor transcript through a
20 hierarchical filtering process (Step 11).

21 *Expertise Required*

22 To effectively implement IgSeqR, individuals must be familiar with computational
23 biology and have basic expertise in navigating a Linux command-line environment.
24 Users will need to be comfortable installing the necessary bioinformatics tools
25 involved, preferably via the conda package manager. The protocol provides annotated
26 scripts to run the pipeline, although proficiency in scripting languages, particularly
27 BASH, and large-scale sequencing data and their data formats (**Table 3**) is beneficial.
28 Familiarity with the principles of immunogenetics, BCR structure and function, and B-
29 cell biology in health and disease is expected for the interpretation and curation of the
30 results (<https://www.imgt.org/IMGTEducation/>). While the protocol can be performed
31 by a skilled graduate student or postdoctoral researcher with the necessary
32 computational expertise, collaboration with a specialized core facility for sequencing
33 analysis may be advantageous when generating and processing primary high-
34 throughput sequencing data.

35 *Limitations of Method*

36 Our initial use of IgSeqR with RNA-seq data from a cohort with unknown tumor B cell
37 percentage [20, 21] demonstrated the utility of our protocol [14]. We used selection
38 criteria that were designed to have the maximal confidence that the sequence
39 identified was tumor-derived (at least 5-fold higher frequency than any other functional
40 full transcript with different CDR3 identified). The full tumor *IGHV-IGHD-IGHJ*
41 sequences including the *IGHC* constant region isotype were defined in 339/489 (69%)
42 samples with RNA-seq data available [14]. However, the probability of identifying the
43 tumor sequence could be maximized by changing certain parameters, including the
44 fold increase of the dominant to the other sequences' frequency or the length of the
45 transcript desired.

1 Although the success rate was lower compared to those with $\geq 10\%$ estimated B cells,
2 a full IG rearrangement could be identified in many samples with $< 10\%$.

3 The main limitation of IgSeqR is accessibility to high-quality RNA and the experimental
4 costs of RNA-seq. However, the costs might be a limitation for large-scale cohorts,
5 and not for well-selected samples. The poor quality of the RNA-seq data is a limitation.
6 RNA extracted from formalin-fixed paraffin-embedded (FFPE) tumor samples, which
7 are commonly available in diagnostic settings, is often of low quality [47], and is
8 currently inadequate for IgSeqR.

9 The sequencing chemistry employed during data generation can influence the outputs
10 of the protocol. IgSeqR protocol has been designed and tested using paired-end
11 sequencing, which is recommended for *de novo* assembly of RNA libraries generated
12 from a polyA library prep and allows the recovery of unmapped reads [48]. The use of
13 sequencing assays and analytical pipelines that remove unmapped reads will severely
14 limit IgSeqR reliability and should therefore not be used.

15 A benchmarking comparison of 10 DLBCL samples demonstrated notably longer
16 runtimes when compared to our CLL cohort, with average runtimes taking 247 minutes
17 in DLBCL vs 33 minutes in CLL per sample (**Table S4**). The cellular complexity and
18 lower tumor purity (**Table S1**) of a DLBCL tissue sample may contribute to these longer
19 runtimes compared to CLL blood samples. However, this is likely to have been
20 compounded by the higher number of starting reads in DLBCL cases (121.2 million on
21 average) compared to CLL (71.1 million on average), which increases the processing
22 requirements at each stage of the protocol.

23 Overall, sample characteristics, sequencing chemistry, and data quality may limit the
24 efficacy of IgSeqR. Quality control assessments should be performed, and any
25 necessary errors should be corrected before using IgSeqR.

26

27 **Materials**

28 *Hardware*

29 The IgSeqR protocol is designed to be versatile, allowing compatibility with various
30 computing resources, ranging from laptops to high-performance computing clusters,
31 and cloud computing platforms. All analyses, including those for comparison with
32 MixCR and TRUST4, were conducted using the Iridis5 high-performance computing
33 cluster at the University of Southampton, utilizing 8 x 2.0 GHz CPU cores and 32 GB
34 RAM to simulate a typical desktop workstation. Default settings were used for MixCR
35 and TRUST4 following the RNA-seq from raw FASTQ files protocols from each tool's
36 documentation.

37 However, the protocol can be run on less powerful hardware with longer expected
38 runtimes. Before starting the protocol, users should carefully consider the exact
39 resources available on their machine, including CPU cores and RAM (considering the
40 RAM utilized by the operating system), to mitigate errors.

41 *Software*

- 42 • Operating system: Linux distribution (tested on Red Hat Enterprise v 7.9 and
43 Ubuntu versions 16, 22 and 24 distributions)

- Conda package manager (<https://conda.io>) to install the IgSeqR environment. All dependencies of IgSeqR are documented in the environment file (**Supplement 1**), which eliminates the need for manual installation of individual tools and dependencies. The main software tools used in IgSeqR are listed below along with their versions as documented in the environment file:
 - BLAST (v 2.13.0) [44]
 - HISAT2 (v 2.2.1) [42]
 - Kallisto (v 0.48.0) [45]
 - Samtools (v 1.16.1) [49]
 - Trinity (v 2.13.2) [43]

To create a conda environment from the command line, navigate to the directory containing the environment file and run the following command:

```
$ conda env create -f environment.yml
```

Replacing 'environment.yml' with the filepath of the environment file.

Once the environment is created, it can be activated by running the following command:

```
$ conda activate IgSeqR
```

- The protocol below provides a detailed explanation of each command required for the operation of the IgSeqR protocol. Each command can be run independently; however, the protocol is designed to be run as a complete pipeline from a Linux shell script. An example BASH script has been provided (**Supplement 2**) which will carry out all analytical steps, if a conda environment containing the necessary software (described above) is correctly setup and the correct experimental variables have been included in the accompanying configuration file (**Supplement 3**). This can be performed by calling the following command in the directory outputs and intermediate files should be written to:

```
$ bash path/to/IgSeqR/Script.sh
```

Where 'path/to/IgSeqR/Script.sh' specifies the location of the IgSeqR script file (**Supplement 2**)

Users must read the protocol thoroughly before performing analysis using the provided scripts to facilitate error debugging and configuration for experiment-specific requirements.

Data

In order to implement this protocol users will need:

- Paired-end RNA sequencing data in either FASTQ or BAM format
- Indexed reference transcriptome for HISAT2 alignment. The protocol was designed and tested using the HISAT2 pre-indexed GRCh38 reference which can be downloaded from the HISAT2 Repository using the command:

```
1 $ wget https://genome-  
2 idx.s3.amazonaws.com/hisat/grch38_snptran.tar.gz
```

3 Alternatively, custom indexed reference from a user provided reference
4 transcriptome can be generated using the `hisat2-build` command, as
5 described in the HISAT2 documentation
6 (<https://daehwankimlab.github.io/hisat2/manual/>)

- 7 • Genomic coordinates associated with target regions for *de novo* transcript
8 assembly. In this application, we have focused on *IG* heavy and light chains
9 coordinates which are supplied in the Procedure section below.
- 10 • Reference sequences for *IG* heavy (*IGHV*, *IGHD*, *IGHJ*) and light (*IGKV*, *IGKJ*,
11 *IGLV*, *IGLJ*) chain genes. The references used to develop this protocol can be
12 found in the supplementary material (**Supplements 4-5**) However, the IMGT
13 database is regularly updated online. Therefore the individual gene reference
14 files should be downloaded from IMGT (**Table 4**) and merged into reference
15 FASTA files for *IG* heavy and *IG* light chains before use of the pipeline.

16
17

18 Procedure

19 *Data pre-processing of newly generated sequencing data (Pre-pipeline)*

20 **1.1** The pipeline has been optimized on fastq files from Illumina sequencing platforms
21 (Illumina, Hayward, CA, USA). Users with newly generated sequencing data in
22 BCL format should follow illumina protocols for converting data into fastq format.
23 Users with fastq files should commence the pipeline at “Step 2. *Genome*
24 *Alignment*”

25 *Pre-processing published and existing sequencing data (~ 5 minutes)*

26 **1.2** FASTQ files are required for downstream steps in this pipeline, however published
27 RNA-seq datasets often provide aligned or unaligned BAM files, in which case
28 FASTQ records must first be extracted from these files, using the `fastq` command
29 from Samtools.

30
31 **CRITICAL STEP:** The `fastq` command requires BAMs to first be sorted by name
32 rather than the default sorting by chromosomal coordinates to ensure proper read
33 pairing. This can be achieved by running the Samtools `sort` command.

34
35 The following example command could be used to sort, and extract FASTQ records
36 from a paired-end BAM file ‘`sample.bam`’. This command uses 8 CPU threads for
37 parallelization, and outputs compressed FASTQ files for read 1, read 2, and
38 unpaired singleton reads to ‘`read1.raw.fastq.gz`’, ‘`read2.raw.fastq.gz`’,
39 respectively.

40

```
1 $ samtools sort -n -@ 8 sample.bam -o sorted.bam
2
3 $ samtools fastq -@ 8 -n -c 6 sorted.bam \
4 -1 read1.raw.fastq.gz \
5 -2 read2.raw.fastq.gz \
6 -0 /dev/null -s /dev/null
```

7
8 The `-n` parameter in `sort` is used to sort BAM file by name, `-@` parameter
9 specifies the number of CPU threads to be used for parallelization of tasks. The `-n`
10 option in `fastq` is used to leave the read names as they are provided. The `-c`
11 option sets the compression level of the output files. 'sample.bam' specifies the
12 path to the input BAM file. `-1` and `-2` specify the desired paths for the compressed
13 FASTQ output files for read 1 and read 2, respectively. `-0 /dev/null` and
14 `-s /dev/null` discard any discarding singletons, supplementary and
15 secondary reads.

16
17 **CRITICAL STEP:** It is important to perform quality control (QC) to ensure that the
18 data is of sufficient quality for downstream analysis. A widely used QC tool is
19 FastQC, which produces a detailed report of several quality metrics including per
20 base sequence quality, per sequence quality scores, per base sequence content,
21 per sequence GC content, and sequence length distribution, among others detailed
22 at in the FastQC documentation [41]. If any issues are identified, corrective
23 measures should be taken as per local procedures or general best practice [50]. If
24 the tumor purity is unknown, it is advisable to estimate this through identification of
25 the B cell proportion using a computational cellular deconvolution tool such as
26 Cibersort [31].

27 *Genome Alignment (~ 10 minutes)*

28 **2.** FASTQ reads are aligned to a reference genome using HISAT2 which produces
29 SAM output file which is processed by Samtools. These commands can be run as
30 a pipeline to save computational resources. The HISAT2 SAM can be passed to
31 Samtools `view` for conversion to BAM format which is then sorted using the
32 Samtools `sort` command. Upon completion of Samtools `sort`, Samtools `index`
33 is run to create an accompanying index file for the BAM.

34
35 The following example command can be used to align FASTQ input files
36 'read1.raw.fastq.gz' and 'read2.raw.fastq.gz' to the GRCh38 reference
37 transcriptome using 8 CPU threads. The resulting HISAT2 aligned BAM file is
38 output as 'hisat_output.bam' and its corresponding index as
39 'hisat_output.bam.bai':

```
40
41 $ hisat2 -p 8 --phred33 -x grch38_snp_tran \
42 -1 read1.raw.fastq.gz -2 read1.raw.fastq.gz | \
43 samtools view -@8 -bS -0 - - | \
44 samtools sort -@8 - -o hisat_output.bam &&
45 samtools index -@8 hisat_output.bam -o hisat_output.bam.bai
```

46
47 The `-p` or `-@` parameters specifies the number of CPU threads to be used for
48 parallelization, while `--phred33` specifies the encoding format of the quality

1 scores. The `-x` parameter specifies the path and basename of the indexed
2 reference transcriptome files. The input FASTQ file paths are specified by `-1` and
3 `-2` for read 1 and read 2, respectively. The SAM is converted to bam using `-bS`
4 with `-0` specifying no additional filtering or format conversions, and `-` signifying
5 the standard input from the previous command. The sorted output is written to a
6 file path specified by `-o hisat_output.bam` from which an index file is created
7 and written to the file path specified by `-o hisat_output.bam.bai`.

8

9 *Read selection (<5 minutes)*

10 **3.** Samtools is used to remove all reads except those that map to the *IG*-associated
11 loci and those that are unmapped from the HISAT2 aligned BAM file, ensuring that
12 highly variable *IG* regions that are difficult to map are retained.

13

14 **TROUBLESHOOTING:** If working from existing data, unmapped reads may have
15 been removed and the alignment files may not contain sufficient *IG* reads to
16 produce quality results.

17

18 The `view` command is used to extract the *IG*-associated loci and unmapped reads
19 independently, before joining using the `merge` command.

20

21 The following example command can be used to filter the HISAT2 aligned BAM file
22 'hisat2_output.bam', retaining reads mapping to the *IGH*, *IGK*, and *IGL* loci
23 and unmapped reads, using 8 threads for parallelization. The resulting filtered bam
24 BAM file is output as 'IG_filtered.bam'. Process substitution can be applied
25 when using a supported Unix shell to avoid the generation of temporary files:

26

```
27 $ samtools merge -f IGH_filtered.bam \  
28 <(samtools view -@ 8 -b -f 4 hisat2_output.bam) \  
29 <(samtools view -@ 8 -b hisat2_output.bam 14:105550000-  
30 106900000 2:87000000-92000000 22:20500000-24500000)
```

31

32 Where `-@` specifies the number of CPU threads to be utilized for parallelization,
33 `-b` specifies the output format as BAM, `-f 4` returns sequences which have the
34 unmapped Samtools flag, `hisat2_output.bam` is the full input HISAT2 aligned
35 BAM file. The *IG* coordinates `14:105550000-106900000`, `2:87000000-92000000`
36 and `22:20500000-24500000` for *IGH*, *IGK* and *IGL*, respectively,
37 are specified in the format `chr:start-end` where `chr` is the chromosome
38 number, `start` is the numerical position of the first nucleotide in the loci and `end`
39 is the numerical position of the last nucleotide.

40 **TROUBLESHOOTING:** The format of the *IG* coordinates will depend on the
41 reference transcriptome used to generate the aligned BAM file. The HISAT2
42 indexed GRCh38 reference uses numerical values for chromosome (e.g., 14).
43 However, other references may also include a 'chr' prefix (e.g., chr14). Additionally,
44 if a different reference transcriptome build is used (e.g., GRCh37) the coordinates
45 should be converted accordingly.

46

1 *De Novo Transcript Assembly (~ 15 minutes)*

2 **4.** Trinity accepts FASTQ input files which must be extracted from the
3 'IGH_filtered.bam' BAM file using the Samtools `fastq` command (as
4 described in Step 1.2).

5
6 The following example command can be used to perform Trinity *de novo* assembly
7 with the input filtered FASTQ files 'IG_filtered_read1.fastq' and
8 'IG_filtered_read2.fastq', using 8 threads for parallelization and 32Gb
9 RAM. The resulting transcriptome FASTA file is output as
10 'trinity_transcripts.fasta':

```
11  
12 $ Trinity -CPU 8 -max_memory 32G -seqType fq \  
13 --left IG_filtered_read1.fastq \  
14 --right IG_filtered_read2.fastq \  
15 --output trinity_transcripts \  
16 --no_normalize_reads \  
17 --min_contig_length 500 \  
18 --full_cleanup
```

19
20 Where `--CPU` specifies the number of CPU threads to be utilized for
21 parallelization, `--max_memory` specifies the maximum memory to be utilized, `--`
22 `seqType fq` specifies that the input files are in FASTQ format, `--left` and `-`
23 `right` are the filtered input FASTQ files for read 1 and read 2, respectively, and
24 `-output <output>` is the basename of the output files.

25
26 **CRITICAL STEP:** Read normalization aims to reduce bias in assembly by down
27 sampling highly expressed reads. Input data will be enriched for *IG* transcripts. This
28 can lead to a reduction of reads for low-abundance transcripts, which can lead to
29 incomplete assembly or loss of rare transcripts and should be disabled using `-`
30 `no_normalize_reads`.

31
32 **CRITICAL STEP:** Short contigs may represent partial or fragmented *IG* transcripts,
33 which can affect downstream analysis and interpretation. Using a minimum contig
34 length of 500 with `-min_contig_length`, most assembled transcripts will
35 contain the full *IGV-(IGD)-IGJ* recombination.

37 *IG Transcript Selection (< 5 mins)*

38 The protocol permits the detection and quantification of *IG* heavy and/or light chains.
39 The steps below provide examples of *IG* heavy chain transcript extraction, but can be
40 adapted to extract the *IG* light chain transcript.

41 **5.** To extract putative *IG* sequences from the Trinity assembly, the transcriptome
42 FASTA file containing the assembled contigs are searched against a reference
43 sequence using BLAST.

44

1 Individual BLAST databases should be generated using the reference sequences
2 for *IG* heavy (*IGHV,IGHD,IGHJ*) (**Supplement 4**) and light (*IGKV,IGKJ,IGLV,IGLJ*)
3 (**Supplement 5**) chains as required using the `makeblastdb` command.

4 The following example command describes how to generate a BLAST database
5 from the *IG* heavy reference FASTA sequences 'IGH_reference.fasta':

```
6  
7 $ makeblastdb -in IGH_reference.fasta -parse_seqids -dbtype  
8 nucl
```

9
10 Where `-in` specifies the input FASTA file containing reference sequences, `-`
11 `parse_seqids` allows the FASTA headers to be parsed along with their
12 sequence, and `-dbtype nucl` specifies the sequence content to be nucleotides.

13
14 **TROUBLESHOOTING:** Ensure that the sequence headers in the FASTA file do
15 not contain the pipe ("|") character as it is a reserved character for the ID parser,
16 which can cause an error.

17
18 **6.** The assembled transcripts are compared against the reference database(s)
19 generated in step 5 using the `BLASTN` command. This produces a tabular output
20 that can be passed to the `cut` and `uniq` commands to obtain a unique list of *IG*
21 transcript IDs, which are used by `samtools faidx` to extract the corresponding
22 sequences from the assembled transcripts FASTA file.

23
24 The following example command can be used to select *IG* transcripts covering
25 reference *IG* heavy FASTA sequences in the 'IGH_reference.fasta' file from
26 the assembled transcripts 'trinity_transcripts.fasta', to produce the
27 filtered FASTA file 'IGH_transcripts.fasta':

```
28  
29 $ blastn -db IGH_reference.fasta \  
30 -query trinity_transcripts.fasta -outfmt 6 | \  
31 cut -f1 | uniq | xargs -n 1 samtools faidx  
32 trinity_transcripts.fasta > IGH_transcripts.fasta
```

33
34 Where, `-db` specifies the path to the FASTA file used to generate the reference
35 database for either *IG* heavy or light sequences, `-query` specifies the path to the
36 FASTA file containing the Trinity assembled transcripts, `-outfmt 6` sets the
37 output format to be tabular, `cut -f1` selects the transcript ID (first) column in the
38 tabular `BLASTN` output, `uniq` removes duplicate transcript IDs, `xargs -n 1`
39 reads the IDs from output of the `uniq` (one ID per line) and passes them to
40 `samtools faidx` as separate arguments.

41 *Transcript Quantification (< 5 minutes)*

42 **7.** Abundance of selected transcripts is quantified using the Kallisto pseudoalignment
43 tool which first requires a Kallisto index to be built from the input FASTA file using
44 the `index` command.

45

1 The following example command can be used to generate a Kallisto index file
2 'kallisto.index' for the *IG* heavy chain filtered transcript FASTA sequence file
3 'IGH_transcripts.fasta':

```
4 kallisto index -i kallisto.index IGH_transcripts.fasta
```

6
7 Where `-i` specifies the filename of the Kallisto index to be constructed and
8 'IGH_transcripts.fasta' is the path to the filtered *IG* transcripts FASTA
9 sequences.

10
11 **8.** The generated index is used in the `quant` command, along with FASTQ files used
12 to assemble the transcripts to quantify the abundance of the *IG* filtered transcripts.

13
14 The following example command can be used to quantify the abundance of
15 transcripts in the *IG* filtered transcript FASTQ files (generated in step 4)
16 'IG_filtered_read1.fastq' and 'IG_filtered_read2.fastq', using 8
17 threads:

```
18 kallisto quant -i kallisto.index -t 8 \  
19 IG_filtered_read1.fastq IG_filtered_read2.fastq
```

21
22 Where `-i` specifies the filename of the Kallisto index, `-t` specifies the number
23 of CPU threads to be utilized for parallelization, and
24 'IG_filtered_read1.fastq' and 'IG_filtered_read2.fastq' are the *IG*
25 filtered FASTQ files for read 1 and read 2, respectively.

26
27 **9.** The five most abundant transcripts IDs are identified based on their transcript per
28 million (TPM) value by passing the Kallisto output through the `tail`, `sort`, `head`
29 and `cut` commands, and their corresponding FASTA sequences are extracted
30 using `samtools faidx` command.

31
32 The following example command can be used to identify the five most abundant
33 transcript IDs from the Kallisto output 'abundance.tsv', extract their
34 corresponding transcript sequences from 'IGH_transcripts.fasta' and write
35 to an output FASTA file called 'IGH_TPM_filtered.fasta':

```
36 $ tail -n +2 abundance.tsv | \  
37 sort -t $'\t' -k5,5nr | head -5 | cut -f1 | \  
38 xargs -n 1 samtools faidx IGH_transcripts.fasta >  
39 IGH_TPM_filtered.fasta
```

41
42 Where `-n +2` selects all rows except the first (header) from the Kallisto
43 quantification output, `-t $'\t'` specifies the delimiter of the input as tab, `-`
44 `k5,5nr` sorts the remaining lines by the fifth column (TPM) in reverse numerical
45 order, `head -5` outputs the first 5 lines of the sorted file and `cut -f1` extracts
46 the first column (IDs) from the output. The IDs are read (one ID per line) using
47 `xargs -n 1` which then passes them to `samtools faidx` as separate
48 arguments.

1 TROUBLESHOOTING: The number of most abundant transcripts to take forward
2 has been suggested as 5. This has been found to strike a good balance between
3 analytical efficiency and identification of the dominant tumor transcript. In instaces
4 where no full-length, productive transcripts are not obtained within the top 5
5 transcripts, users may wish to increase the number of transcripts to take forward
6 for analysis.

7
8

9 *Dominant IG Transcript selection (~ 15 minutes)*

10

11 **10.** The top 5 most abundant transcripts identified in step 9 will be submitted to the
12 IMGT/V-QUEST tool (https://imgt.org/IMGT_vquest/input) for sequence analysis
13 and annotation. In the sequence submission section of the IMGT/V-QUEST tool,
14 the top 5 transcript sequences should be provided either by copy and pasting the
15 sequences from the FASTA file or by directly uploading the FASTA file. The
16 parameters 'Species' and 'Receptor type or locus' should be set to 'Homo sapiens
17 (human)' and 'IG', respectively. Finally, the output format should be set to 'C.Excel
18 file'. The IMGT/V-QUEST tool will annotate and analyze the submitted sequences
19 for their corresponding *IGV*, *IGHD* (for the heavy chain only) and *IGJ* genes, their
20 junction at the *CDR3* region, and other related features.

21

22 **11.** The outputs of the Kallisto quantification and IMGT/V-QUEST results transcript
23 are used to identify the dominant/consensus (Tumor) *IG* transcript present within
24 the RNA-seq dataset. This process may require manual interpretation but follows
25 the following hierarchical filtering criteria:

- 26 i. Presence of a full transcript sequence (from codon 1 in FR1 to codon
27 129 in FR4 included), identified by IMGT/VQUEST.
- 28 ii. Presence of 'productive' V-domain functionality call by IMGT/V-QUEST
- 29 iii. The highest estimated read count (est. count) determined by Kallisto.
- 30 iv. The est. count is greater than 5-fold higher than any of the other 4
31 transcripts selected if different. A reduction of the fold amount difference
32 will increase the probability to identify a “dominant” sequence in cases
33 with low tumor infiltration.
- 34 v. The ability to determine the *IG* constant region class and subclass.

35 **Timing**

36 Benchmarking was conducted using the computational hardware described in the
37 materials section. The dominant *IG* heavy and light chain transcripts were extracted
38 from FASTQ files generated from high-purity CLL samples with an average starting
39 read count of 71.1 million following initial HISAT2 alignment. In similar conditions, the
40 full pipeline can be expected to take less than 1 hour per sample. Specific timings can
41 be found in the procedure section headers for each stage of the analytical pipeline.
42 The duration of each stage may vary depending on the input file type (BAM files require
43 additional pre-processing), hardware used to run the pipeline, heterogeneity of B-cell
44 populations, and number of starting sequencing reads generated from the samples.

45 **Anticipated results**

1 Upon successful completion of the IgSeqR protocol, users will have generated the
2 following output files for *IG* heavy and/or light chain transcripts:

- 3 • The five most abundant assembled *IG* transcripts in FASTA format
- 4 • Table of quantifications for these *IG* transcripts in tsv format
- 5 • Annotations for the top five *IG* transcripts generated by IMGT/V-QUEST.

6 For further insights, users can refer to our previously published [14], which includes
7 results and examples of downstream analysis.

8 **Future Applications**

9 Future work is planned to develop further the existing protocol and evaluate its efficacy
10 for deriving smaller, less dominant, clonal populations to widen the application of the
11 protocol. When the tumor *IG* sequence is already known, we will apply this approach
12 for the determination of the minimal residual disease in repeat samples following anti-
13 cancer therapy.

14 We will also investigate the protocol's potential use with RNA-seq data generated from
15 FFPE material. However, there are intrinsic limitations of RNA-seq data quality from
16 FFPE, and areas where optimization or adaptation may be necessary will need to be
17 identified.

18 Future work will also focus on the annotation refinement of the *IGC* region. This work
19 will facilitate and accompany the development of the protocol into a comprehensive
20 bioinformatics tool for immunobiologists.

21 The protocol will also be investigated for its use in any other genomic regions that are
22 challenging to map to a reference genome, including the T-cell Receptor specific
23 fusion or deregulating gene rearrangements that are not represented in the reference
24 transcriptome [51].

25

26 **Supplementary information**

- 27 • Supplementary Tables.xlsx
- 28 • Supplement 1. IgSeqR Environment.yml
- 29 • Supplement 2. IgSeqR BASH Script.sh
- 30 • Supplement 3. IgSeqR Configuration File
- 31 • Supplement 4. IGH References.fasta
- 32 • Supplement 5. IGKL References.fasta

33

34

35 **Author contributions statements**

36 D.B. and B.J.S. designed IgSeqR bioinformatic pipeline, analyzed and interpreted
37 data, and wrote the manuscript. D.T. and G.C analyzed, interpreted data, and
38 contributed to the immunoglobulin gene analysis pipeline validation. B.S. , A.O. and
39 J.B contributed to the analysis and interpretation of the data. F.F. designed the study,
40 supervised research, interpreted data and wrote the manuscript. All authors reviewed
41 and approved the manuscript.

1

2 ORCID for corresponding authors:

3 Francesco Forconi <https://orcid.org/0000-0002-2211-1831>

4

5 Acknowledgments

6 The authors are grateful to the Faculty of Medicine Tissue Bank (Cancer Sciences,
7 University of Southampton) for the processing and storage of the primary lymphoma
8 specimens. This work was supported by Cancer Research UK (ECRIN-M3 accelerator
9 award C42023/A29370, and BTERP project C36811/A29101). D.T. was funded by the
10 Eyles Cancer Immunology PhD scholarship. G.C. was funded by the Eyles Cancer
11 Immunology Fellowship and the Southampton Cancer Immunology Centre Pump-
12 priming award 2021). Genetic data for IgSeqR protocol were obtained via the National
13 Cancer Institute Genomic Data Commons for Genotypes and Phenotypes (accession
14 phs001444.v1.p1). The authors acknowledge the use of the IRIDIS High Performance
15 Computing Facility, and associated support services at the University of Southampton,
16 in the completion of this work.

17 Competing interests

18 The authors declare no potential conflicts of interest.

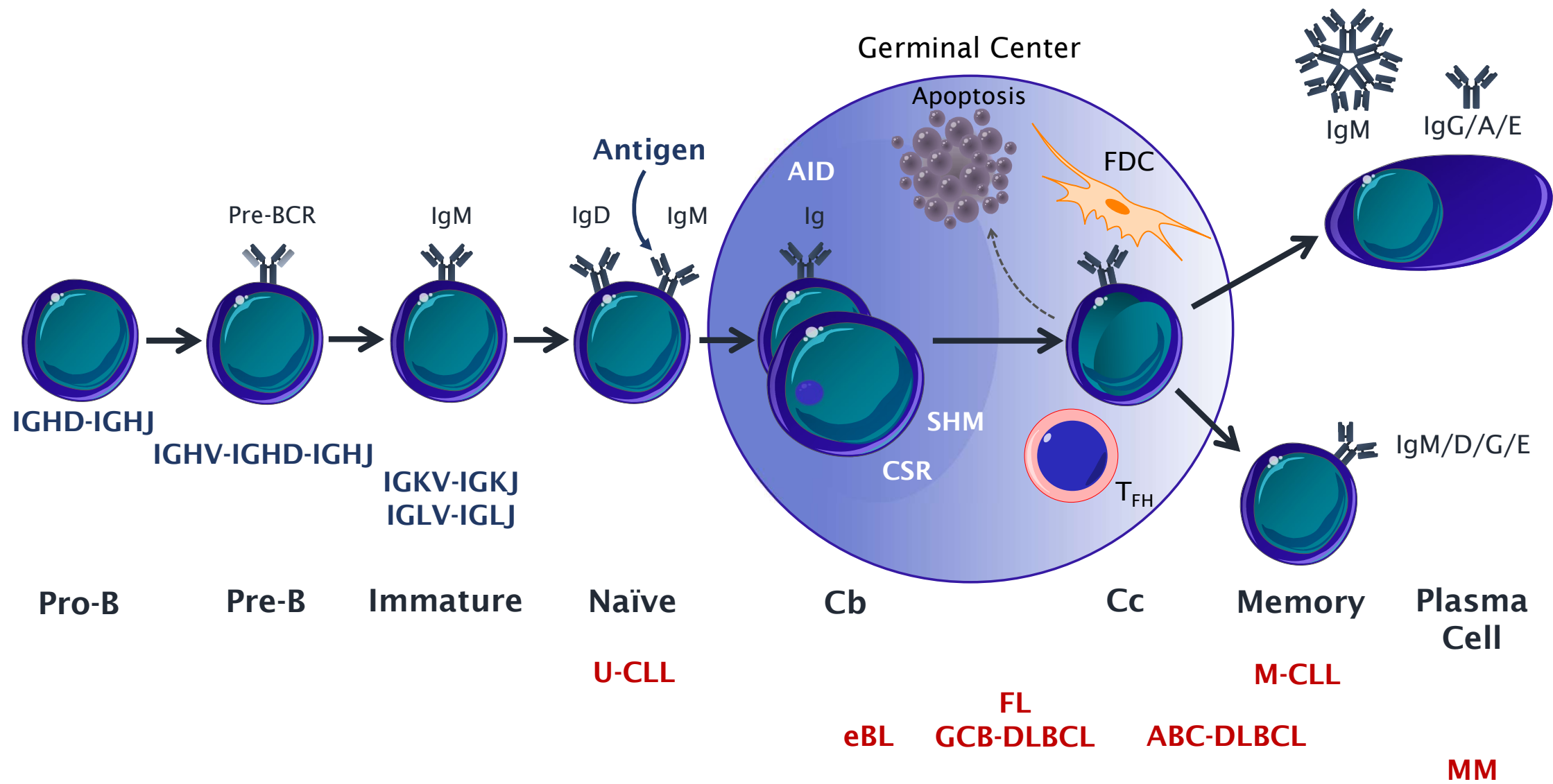
19

20 References

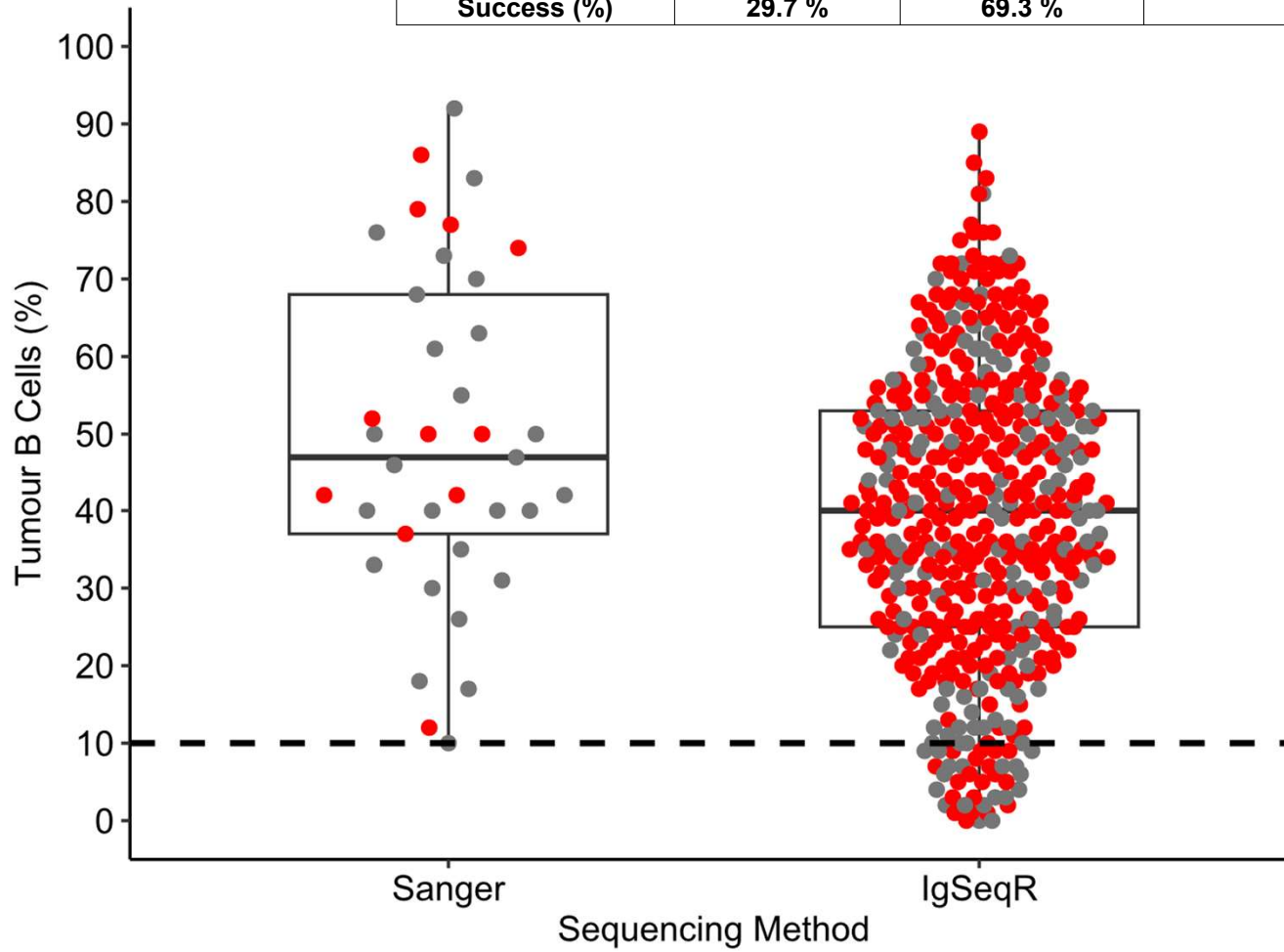
- 21 1. Lam, K.P., R. Kuhn, and K. Rajewsky, *In vivo ablation of surface immunoglobulin on*
22 *mature B cells by inducible gene targeting results in rapid cell death.* *Cell*, 1997. **90**(6):
23 p. 1073-83.
- 24 2. Stevenson, F.K., et al., *The occurrence and significance of V gene mutations in B cell-*
25 *derived human malignancy.* *Adv Cancer Res*, 2001. **83**: p. 81-116.
- 26 3. Vitora, G.D. and M.C. Nussenzweig, *Germinal Centers.* *Annual Review of*
27 *Immunology*, 2022. **40**(1): p. 413-442.
- 28 4. Forconi, F., S.A. Lanham, and G. Chiodin, *Biological and Clinical Insight from Analysis*
29 *of the Tumor B-Cell Receptor Structure and Function in Chronic Lymphocytic*
30 *Leukemia.* *Cancers (Basel)*, 2022. **14**(3).
- 31 5. Stevenson, F.K., F. Forconi, and T.J. Kipps, *Exploring the pathways to chronic*
32 *lymphocytic leukemia.* *Blood*, 2021. **138**(10): p. 827-835.
- 33 6. Efremov, D.G., S. Turkalj, and L. Laurenti, *Mechanisms of B Cell Receptor Activation*
34 *and Responses to B Cell Receptor Inhibitors in B Cell Malignancies.* *Cancers*, 2020.
35 **12**(6): p. 1396.
- 36 7. Forconi, F., et al., *The normal IGHV1-69-derived B-cell repertoire contains stereotypic*
37 *patterns characteristic of unmutated CLL.* *Blood*, 2010. **115**(1): p. 71-7.
- 38 8. Seifert, M., et al., *Cellular origin and pathophysiology of chronic lymphocytic leukemia.*
39 *J Exp Med*, 2012. **209**(12): p. 2183-98.
- 40 9. Damle, R.N., et al., *Ig V gene mutation status and CD38 expression as novel*
41 *prognostic indicators in chronic lymphocytic leukemia.* *Blood*, 1999. **94**(6): p. 1840-7.
- 42 10. Hamblin, T.J., et al., *Unmutated Ig V(H) genes are associated with a more aggressive*
43 *form of chronic lymphocytic leukemia.* *Blood*, 1999. **94**(6): p. 1848-54.
- 44 11. Niemann, C.U., et al., *Fixed-duration ibrutinib–venetoclax versus chlorambucil–*
45 *obinutuzumab in previously untreated chronic lymphocytic leukaemia (GLOW): 4-year*
46 *follow-up from a multicentre, open-label, randomised, phase 3 trial.* *The Lancet*
47 *Oncology*, 2023.

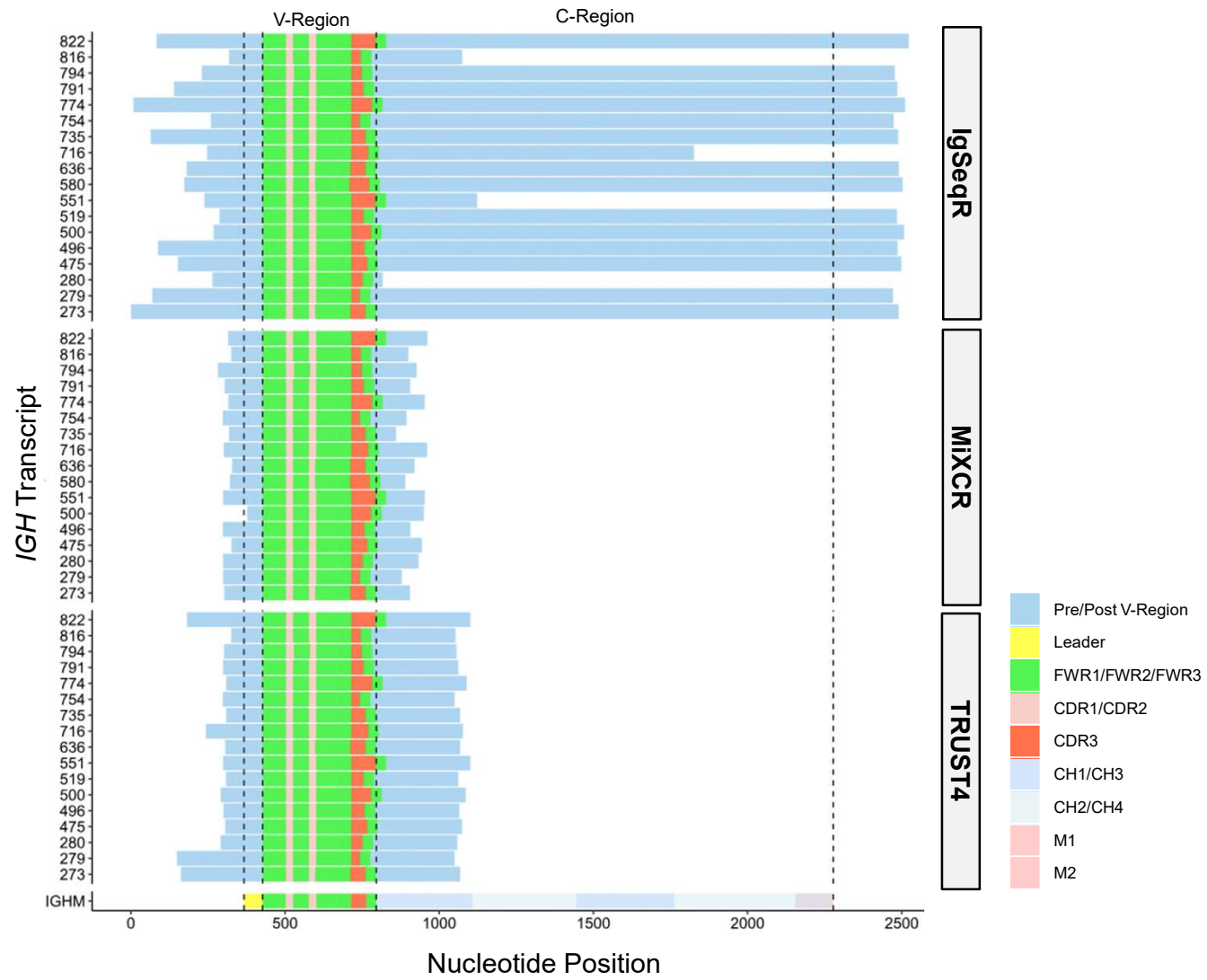
- 1 12. Stevenson, F.K. and F. Forconi, *The essential microenvironmental role of*
2 *oligomannoses inserted into the antigen-binding sites of lymphoma cells*. Blood, 2023.
- 3 13. Zhu, D., et al., *Acquisition of potential N-glycosylation sites in the immunoglobulin*
4 *variable region by somatic mutation is a distinctive feature of follicular lymphoma*.
5 Blood, 2002. **99**(7): p. 2562-2568.
- 6 14. Chiodin, G., et al., *Insertion of atypical glycans into the tumor antigen-binding site*
7 *identifies DLBCLs with distinct origin and behavior*. Blood, 2021. **138**(17): p. 1570-
8 1582.
- 9 15. Coelho, V., et al., *Glycosylation of surface Ig creates a functional bridge between*
10 *human follicular lymphoma and microenvironmental lectins*. Proc Natl Acad Sci U S A,
11 2010. **107**(43): p. 18587-92.
- 12 16. Linley, A., et al., *Lectin binding to surface Ig variable regions provides a universal*
13 *persistent activating signal for follicular lymphoma cells*. Blood, 2015. **126**(16): p. 1902-
14 10.
- 15 17. Odabashian, M., et al., *IGHV sequencing reveals acquired N-glycosylation sites as a*
16 *clonal and stable event during follicular lymphoma evolution*. Blood, 2020. **135**(11): p.
17 834-844.
- 18 18. Sutton, L.A., et al., *Immunoglobulin genes in chronic lymphocytic leukemia: key to*
19 *understanding the disease and improving risk stratification*. Haematologica, 2017.
20 **102**(6): p. 968-971.
- 21 19. Wang, Z., M. Gerstein, and M. Snyder, *RNA-Seq: a revolutionary tool for*
22 *transcriptomics*. Nat Rev Genet, 2009. **10**(1): p. 57-63.
- 23 20. Schmitz, R., et al., *Genetics and pathogenesis of diffuse large B-cell lymphoma*. New
24 England Journal of Medicine, 2018. **378**(15): p. 1396-1407.
- 25 21. Wright, G.W., et al., *A Probabilistic Classification Tool for Genetic Subtypes of Diffuse*
26 *Large B Cell Lymphoma with Therapeutic Implications*. Cancer Cell, 2020. **37**(4): p.
27 551-568.e14.
- 28 22. Bryant, D., et al., *Network analysis reveals a major role for 14q32 cluster miRNAs in*
29 *determining transcriptional differences between IGHV-mutated and unmutated CLL*.
30 Leukemia, 2023. **37**(7): p. 1454-1463.
- 31 23. Teo, Q.W., et al., *Stringent and complex sequence constraints of an IGHV1-69 broadly*
32 *neutralizing antibody to influenza HA stem*. Cell Reports, 2023. **42**(11): p. 113410.
- 33 24. Shiroishi, M., *Structural Basis of a Conventional Recognition Mode of IGHV1-69*
34 *Rheumatoid Factors*, in *Protein Reviews : Volume 21*, M.Z. Atassi, Editor. 2021,
35 Springer International Publishing: Cham. p. 171-182.
- 36 25. Forconi, F., et al., *Insight into the potential for DNA idiotypic fusion vaccines designed*
37 *for patients by analysing xenogeneic anti-idiotypic antibody responses*. Immunology,
38 2002. **107**(1): p. 39-45.
- 39 26. Stevenson, G.T., E.V. Elliott, and F.K. Stevenson, *Idiotypic determinants on the*
40 *surface immunoglobulin of neoplastic lymphocytes: a therapeutic target*. Fed Proc,
41 1977. **36**(9): p. 2268-71.
- 42 27. Hawkins, R.E., et al., *Idiotypic vaccination against human B-cell lymphoma. Rescue of*
43 *variable region gene sequences from biopsy material for assembly as single-chain Fv*
44 *personal vaccines*. Blood, 1994. **83**(11): p. 3279-88.
- 45 28. McCann, K., et al., *Idiotype gene rescue in follicular lymphoma*. Methods Mol Med,
46 2005. **115**: p. 145-71.
- 47 29. Ottensmeier, C.H. and F.K. Stevenson, *Isotype switch variants reveal clonally related*
48 *subpopulations in diffuse large B-cell lymphoma*. Blood, 2000. **96**(7): p. 2550-6.
- 49 30. Ottensmeier, C.H., et al., *Analysis of VH genes in follicular and diffuse lymphoma*
50 *shows ongoing somatic mutation and multiple isotype transcripts in early disease with*
51 *changes during disease progression*. Blood, 1998. **91**(11): p. 4292-9.
- 52 31. Newman, A.M., et al., *Robust enumeration of cell subsets from tissue expression*
53 *profiles*. Nature Methods, 2015. **12**(5): p. 453-457.

- 1 32. Blachly, J.S., et al., *Immunoglobulin transcript sequence and somatic hypermutation*
2 *computation from unselected RNA-seq reads in chronic lymphocytic leukemia*. Proc
3 Natl Acad Sci U S A, 2015. **112**(14): p. 4322-7.
- 4 33. Bolotin, D.A., et al., *MiXCR: software for comprehensive adaptive immunity profiling*.
5 Nat Methods, 2015. **12**(5): p. 380-1.
- 6 34. Canzar, S., et al., *BASIC: BCR assembly from single cells*. Bioinformatics, 2017. **33**(3):
7 p. 425-427.
- 8 35. Kuchenbecker, L., et al., *IMSEQ--a fast and error aware approach to immunogenetic*
9 *sequence analysis*. Bioinformatics, 2015. **31**(18): p. 2963-71.
- 10 36. Mandric, I., et al., *Profiling immunoglobulin repertoires across multiple human tissues*
11 *using RNA sequencing*. Nat Commun, 2020. **11**(1): p. 3126.
- 12 37. Mose, L.E., et al., *Assembly-based inference of B-cell receptor repertoires from short*
13 *read RNA sequencing data with V'DJer*. Bioinformatics, 2016. **32**(24): p. 3729-3734.
- 14 38. Rizzetto, S., et al., *B-cell receptor reconstruction from single-cell RNA-seq with*
15 *VDJPuzzle*. Bioinformatics, 2018. **34**(16): p. 2846-2847.
- 16 39. Song, L., et al., *TRUST4: immune repertoire reconstruction from bulk and single-cell*
17 *RNA-seq data*. Nat Methods, 2021. **18**(6): p. 627-630.
- 18 40. Upadhyay, A.A., et al., *BALDR: a computational pipeline for paired heavy and light*
19 *chain immunoglobulin reconstruction in single-cell RNA-seq data*. Genome Med, 2018.
20 **10**(1): p. 20.
- 21 41. Andrews, S., *FastQC: A Quality Control Tool for High Throughput Sequence Data*.
22 Babraham Bioinformatics, 2010.
- 23 42. Kim, D., et al., *Graph-based genome alignment and genotyping with HISAT2 and*
24 *HISAT-genotype*. Nature Biotechnology, 2019. **37**(8): p. 907-915.
- 25 43. Grabherr, M.G., et al., *Full-length transcriptome assembly from RNA-Seq data without*
26 *a reference genome*. Nat Biotechnol, 2011. **29**(7): p. 644-52.
- 27 44. Altschul, S.F., et al., *Basic local alignment search tool*. J Mol Biol, 1990. **215**(3): p. 403-
28 10.
- 29 45. Bray, N.L., et al., *Near-optimal probabilistic RNA-seq quantification*. Nat Biotechnol,
30 2016. **34**(5): p. 525-7.
- 31 46. Brochet, X., M.P. Lefranc, and V. Giudicelli, *IMGT/V-QUEST: the highly customized*
32 *and integrated system for IG and TR standardized V-J and V-D-J sequence analysis*.
33 Nucleic Acids Res, 2008. **36**(Web Server issue): p. W503-8.
- 34 47. Cazzato, G., et al., *Formalin-Fixed and Paraffin-Embedded Samples for Next*
35 *Generation Sequencing: Problems and Solutions*. Genes (Basel), 2021. **12**(10).
- 36 48. Haas, B.J., et al., *De novo transcript sequence reconstruction from RNA-seq using the*
37 *Trinity platform for reference generation and analysis*. Nat Protoc, 2013. **8**(8): p. 1494-
38 512.
- 39 49. Li, H., et al., *The Sequence Alignment/Map format and SAMtools*. Bioinformatics,
40 2009. **25**(16): p. 2078-9.
- 41 50. Hesketh, A.R., *RNA Sequencing Best Practices: Experimental Protocol and Data*
42 *Analysis*, in *Yeast Systems Biology: Methods and Protocols*, S.G. Oliver and J.I.
43 Castrillo, Editors. 2019, Springer New York: New York, NY. p. 113-129.
- 44 51. Haas, B.J., et al., *Accuracy assessment of fusion transcript detection via read-mapping*
45 *and de novo fusion transcript assembly-based methods*. Genome Biology, 2019. **20**(1):
46 p. 213.



	Sanger (n = 37)	IgSeqR (n = 489)	P-value
■ IG Identified	11	339	<0.0001
■ IG Not Identified	26	150	
Success (%)	29.7 %	69.3 %	





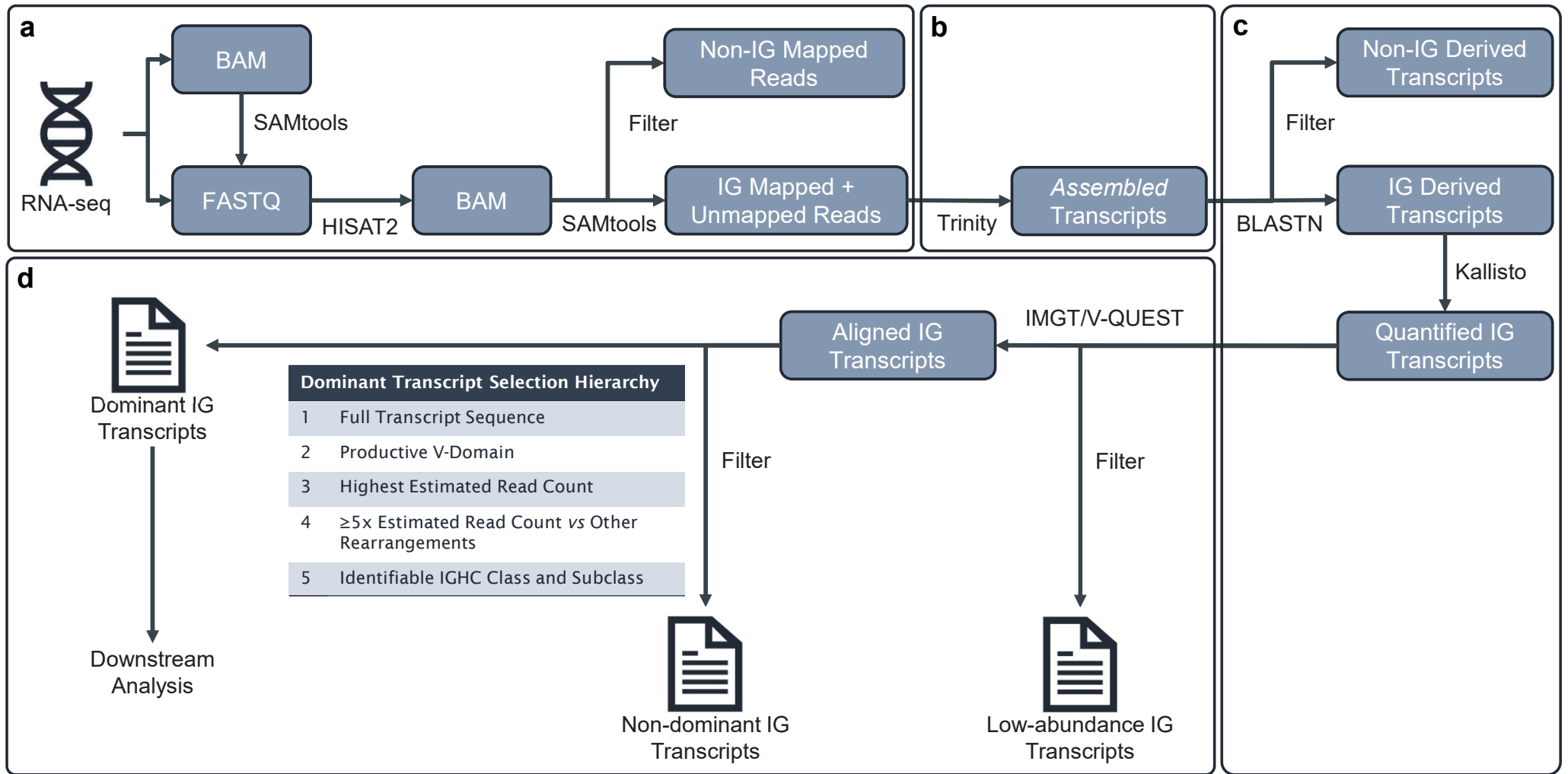


Figure 1. Immunoglobulin gene analysis provides insight into the cell of origin and behavior of B-Cell malignancies.

The immunoglobulin heavy-chain gene repertoire comprises ~51 functional variable (*IGHV*), ~21 diversity (*IGHD*), and ~7 joining (*IGHJ*) genes at the *14q32* locus. In the bone marrow, progenitor B cells (Pro-B cells) undergo an *IGHD-IGHJ* rearrangement. If successful, a complete *IGHV-IGHD-IGHJ* rearrangement occurs at the precursor B cell (Pre-B), which expresses a precursor B-cell receptor (pre-BCR) containing a surrogate VpreB1 light chain and a full IG heavy chain. The pre-BCR will promote rearrangement of *IGKV-IGKJ* at the *2p11.2* locus, and, if this is non-functional in both alleles, the rearrangement of *IGLV-IGLJ* will occur at the *22q11.22* locus in immature B cells. A successful rearrangement of the IG light chain enables the expression of a competent immunoglobulin M (IgM) and autoreactive B cell clones within the bone marrow microenvironment will be deleted, ensuring the production of functional non-autoreactive naïve B cells expressing IgM and IgD. IgM+ve IgD+ve naïve B cells exit the bone marrow and migrate to peripheral lymphoid organs (spleen, lymph nodes, MALTs, etc.) where they will encounter antigen, and they will undergo class-switch recombination (CSR) and somatic hypermutation (SHM) in the presence of activation-induced cytidine-deaminase (AID) in a germinal center (GC) reaction at the centroblast (Cb) stage (dark zone). During SHM, Cb introduce point mutations in the IG variable region genes to mature affinity to antigen. Centrocytes (Cc) emerge in the light zone where their fate will depend on their BCR interactions with immune complexes on follicular dendritic cells (FDC) in the presence of T follicular helper (T_{FH}) cells. Cc with the BCR of the right affinity to antigen receive survival signals and differentiate into memory B cells or plasma cells, while the others will undergo apoptosis.

The tumor *IG* genes preserve the features of the cell having undergone transformation. Chronic Lymphocytic Leukemias with unmutated *IG* genes (U-CLL) arise from pre-GC B-cells and have an aggressive clinical course, while those with mutated *IG* genes (M-CLL) arise from post-GC B cells and display an indolent clinical course. In endemic Burkitt lymphoma (eBL), FL, and some DLBCL, there is intraclonal heterogeneity of the *IGV* gene sequences to indicate that the SHM process is ongoing, as in a GC B cell. Diffuse Large B-cell Lymphoma (DLBCL) can be classified into two major subtypes: GC B-cell-like (GCB) and activated B-cell-like (ABC). Asparagine-x-serine/threonine N-glycosylation motifs (where X is any amino acid except proline) are introduced by SHM, allowing occupation of the sites by oligomannose-type glycans in almost all FL and in ~30% of all GCB-DLBCL. Multiple myeloma (MM) is characterized by the clonal expansion of plasma cells, which carry mutated *IG* and secrete a monoclonal IG in the serum (paraprotein).

Figure 2. Tissue-derived lymphoma samples, in which the tumor *IG* sequence was sought by Sanger or IgSeqR.

Each dot identifies a sample: red dots indicate the samples where the tumor *IG* sequence was identified; grey dots indicate the samples where the tumor *IG* sequence could not be identified. The accompanying table indicates the number and proportion of *IG* sequences identified using the individual methods. There was a significantly higher proportion and probability of identifying the tumor *IG* sequence by RNA-seq/IgSeqR (69%) compared to PCR/Sanger (30%) (X-square with Yates' correction, p-value <0.0001). Sequencing by Sanger was performed only in samples with >10% tumor infiltration by immunophenotype, while IgSeqR was applied to any sample irrespective of (tumor) B cell percentage, as estimated by Cibersort.

Figure 3. Comparison between the transcripts recovered by the IgSeqR pipeline with MiXCR, TRUST4 and a reference *IGHM* transcript.

A direct comparison of three analytical tools for the recovery of *IGHV-IGHD-IGHJ* transcripts recovered from unselected bulk high throughput RNA sequencing data from 18 chronic lymphocytic leukemia samples with high tumor purity. The tools, IgSeqR, MiXCR (v 4.3.2), and TRUST4 (v1.0.12) were run using the Iridis5 high-performance computing cluster at the University of Southampton, utilizing 8 x 2.0 GHz CPU cores and 32 GB RAM to simulate a typical desktop workstation. The resulting transcripts were assessed for recovery of a full-length, productive *IGHV-IGHD-IGHJ* (V-Region) transcript and concordance with matched Sanger sequencing in the V-region. MiXCR recovered *IGHV-IGHD-IGHJ* transcripts for all 18 of the samples, with 17 (94%) having productive and full V-Region coverage, however only 17 used the same IGHV of Sanger, and 14 (78%) had 100% identity with Sanger. TRUST4 generated *IGHV-IGHD-IGHJ* transcripts from 17 (94%) of the samples, all of which had productive and full V-Region coverage and full concordance with Sanger. IgSeqR demonstrated productive and full V-Region coverage and full concordance with Sanger in all 18 (100%) samples. IgSeqR also produced the longest tumor transcripts, averaging a length of 2036 nucleotides, compared to 589 and 769 nucleotides by MiXCR and TRUST4 respectively. Notably, the majority (78%) of the IgSeqR transcripts were long enough to cover the full *IGHM* transcript from leader to the membrane domains (M1 and M2) of the constant region (C-Region), a feature not possible in the shorter transcripts generated by MiXCR or TRUST4.

Figure 4. Schematic representation of the IgSeqR Pipeline.

The experimental design of IgSeqR is divided into four key stages: (a) data pre-processing – RNA sequencing data (RNA-seq) can be supplied in either BAM or FASTQ format. The data are re-aligned to a reference transcriptome by HISAT2, producing a BAM file which is filtered

to retain reads mapping to *IG* gene coordinates, and reads unable to be mapped to the reference; **(b)** *de novo* transcriptome assembly - Trinity is used to *de novo* assemble transcripts from the filtered BAM file; **(c)** *IG* transcript selection and quantification – the assembled transcripts are run through a BLAST query to identify transcripts overlapping *IG* reference sequences. The abundance of the *IG*-derived transcripts is then estimated using Kallisto pseudoalignment; **(d)** *IG* transcript annotation and interpretation – the five most abundant transcripts by TPM are then run through IMG_T/VEQUEST for *IG* alignment and annotation which is used to recover the putative tumor/dominant *IG* transcript using a 5 step hierarchical selection process.

Table 1. Published tools for IG analysis from bulk and single-cell RNA sequencing (RNA-seq)

Tool	Description	Receptor	Sequencing Data	Reference
IG_ID	<i>De Novo</i> assembly of BCR transcripts from bulk RNA-seq data	BCR	Bulk	Blachly et al 2015 [31]
MiXCR	Analysis of raw T- or B- cell receptor repertoire sequencing data	BCR/TCR	Bulk/Single Cell	Bolotin et al 2015 [32]
BASIC	Bayesian inference of immunoglobulin sequences. It offers functionalities for V(D)J gene identification, clonotype analysis, and mutation profiling.	BCR	Single Cell	Canzar et al 2017 [33]
IMSEQ	Provides functionalities for the identification and quantification of IG genes, as well as the detection of somatic hypermutations	BCR/TCR	Bulk	Kuchenbecker et al 2015 [34]
ImReP	Extraction of receptor reads from sequencing data and assemble clonotypes, detect corresponding V(D)J recombinations and correct PCR sequencing errors	BCR/TCR	Bulk	Mandric et al 2020 [35]
V'DJer	Customized read extraction, assembly and V(D)J rearrangement detection and filtering to produce contigs representing the most abundant portions of the BCR repertoire	BCR	Bulk	Mose et al 2016 [36]
VDJPuzzle	Provides a user-friendly interface for the identification of V(D)J rearrangements, clonotype analysis, and visualization of TCR and BCR repertoires.	BCR/TCR	Single Cell	Rizzetto et al 2018 [37]
TRUST4	Performs de novo assembly on V, J, C genes including the hypervariable complementarity-determining region 3 (CDR3) and reports consensus of BCR/TCR sequences	BCR/TCR	Bulk/Single Cell	Song et al 2021 [38]
BALDR	Infers the clonal structure of B-cell repertoires, providing information on clonal abundance, V(D)J gene usage, and somatic hypermutations	BCR	Single Cell	Upadhyay et al 2018 [39]

Table 2. A comparison between RNA-seq based Immunoglobulin Gene analysis tools, IgSeqR MiXCR and TRUST4.

Property	IgSeqR	MixCR	TRUST4
Recognized as IG (IMGT)	18 (100 %)	18 (100 %)	17 (94.44 %)
Productive Sequence	18 (100 %)	18 (100 %)	17 (94.44 %)
Complete VDJ	18 (100 %)	17 (94.44 %)	17 (94.44 %)
IGHV Gene match	18 (100 %)	17 (94.44 %)	17 (94.44 %)
IGHV Seq match	18 (100 %)	14 (77.78 %)	17 (94.44 %)
IGHD Gene Allele match	18 (100 %)	18 (100 %)	17 (94.44 %)
IGHD Seq match	18 (100 %)	18 (100 %)	17 (94.44 %)
IGHJ Gene Allele match	18 (100 %)	18 (100 %)	17 (94.44 %)
IGHJ Seq match	18 (100 %)	18 (100 %)	17 (94.44 %)
CDR3 Seq match	18 (100 %)	18 (100 %)	17 (94.44 %)
Full Sanger VDJ Concordance	18 (100 %)	14 (77.78 %)	17 (94.44 %)
Average Length	2036	589	768
Assembly efficiency (Seconds/Nucleotide)	1.18	8.10	1.44

In red are identified the properties of MixCR or TRUST4 with inferior performance compared to IgSeqR.

Table 3. List of acronyms used in the IgSeqR protocol

Acronyms	Name	Description
BAM	Binary Alignment/Map	The BAM format is a binary representation of sequence alignment data. It is commonly used in genomics to store the results of sequence alignment algorithms.
BASH	Bourne Again Shell	Unix shell and a command language interpreter. It is a default command interpreter on most GNU/Linux systems. Bash can also read and execute commands from a file, called a shell script.
CPU	Central Processing Unit	The CPU is the most important processor in a given computer, responsible for performing basic arithmetic, logic, controlling, and input/output (I/O) operations specified by the instructions in a program
CSV	Comma-Separated Values	CSV is a simple file format used to store tabular data, such as a spreadsheet or database. Each line of the file represents a row of the table, and the values are separated by commas. CSV files are widely supported by spreadsheet and database software, making them easy to import and export data.

FASTA	FASTA Sequence Format	The FASTA format is a text-based format for representing nucleotide or protein sequences. It consists of a single-line description followed by lines of sequence data. The format is widely used in bioinformatics for storing and exchanging sequence data.
FASTQ	FASTQ Sequence Format	The FASTQ format is a text-based format for storing both a biological sequence (usually nucleotide sequence) and its corresponding quality scores. It is widely used to represent raw sequencing data from high-throughput sequencing platforms.
TSV	Tab-Separated Values	TSV is a file format similar to CSV, but with tab characters as the field separator instead of commas. TSV files are commonly used for storing and exchanging tabular data, especially when the data may contain commas or other special characters.
RAM	Random Access Memory	A temporary memory bank in a computer where data which requires quick access is stored. It keeps data easily accessible so a computers processor can quickly find it without having to go into long-term storage to complete immediate processing tasks.

Table 4. Links to the most recent IMGT/V-QUEST reference immunoglobulin heavy and light chain FASTA sequences

Chain	Gene	IMGT Link
Heavy	IGHV	imgt.org/download/V-QUEST/IMGT_V-QUEST_reference_directory/Homo_sapiens/IG/IGHV.fasta
	IGHD	imgt.org/download/V-QUEST/IMGT_V-QUEST_reference_directory/Homo_sapiens/IG/IGHD.fasta
	IGHJ	imgt.org/download/V-QUEST/IMGT_V-QUEST_reference_directory/Homo_sapiens/IG/IGHJ.fasta
Light	IGKV	imgt.org/download/V-QUEST/IMGT_V-QUEST_reference_directory/Homo_sapiens/IG/IGKV.fasta
	IGKJ	imgt.org/download/V-QUEST/IMGT_V-QUEST_reference_directory/Homo_sapiens/IG/IGKJ.fasta
	IGLV	imgt.org/download/V-QUEST/IMGT_V-QUEST_reference_directory/Homo_sapiens/IG/IGLV.fasta
	IGLJ	imgt.org/download/V-QUEST/IMGT_V-QUEST_reference_directory/Homo_sapiens/IG/IGLJ.fasta