



Single-cell Rapid Capture Hybridization sequencing to reliably detect isoform usage and coding mutations in targeted genes

Hongke Peng, Jafar S. Jabbari, Luyi Tian, et al.

Genome Res. published online January 10, 2025

Access the most recent version at doi:[10.1101/gr.279322.124](https://doi.org/10.1101/gr.279322.124)

P<P	Published online January 10, 2025 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

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

Published by Cold Spring Harbor Laboratory Press

1 **Single-cell Rapid Capture Hybridization sequencing reliably detects**
2 **isoform usage and coding mutations in targeted genes**

3

4 **Authors:** Hongke Peng^{1,2}, Jafar S. Jabbari^{1,2}, Luyi Tian^{1,2,8}, Changqing Wang^{1,2}, Yupei You^{1,2}, Chong
5 Chyn Chua^{1,2,3,4}, Natasha S. Anstee^{1,2}, Noorul Amin^{1,2}, Andrew H. Wei^{1,2,5}, Nadia M. Davidson^{1,2},
6 Andrew W. Roberts^{1,2,5}, David C. S. Huang^{1,2}, Matthew E. Ritchie^{1,2,*}, Rachel Thijssen^{1,2,6,7,*}

7

8 ¹*The Walter and Eliza Hall Institute of Medical Research, Melbourne, Australia*

9 ²*Department of Medical Biology, University of Melbourne, Melbourne, Australia*

10 ³*Monash Haematology, Monash Health, Melbourne Australia*

11 ⁴*Clinical Haematology, Northern Health, Melbourne, Australia*

12 ⁵*Department of Clinical Haematology, Royal Melbourne Hospital and Peter MacCallum Cancer
13 Centre, Melbourne, Australia*

14 ⁶*Department of Hematology, Amsterdam UMC, Amsterdam, the Netherlands.*

15 ⁷*Cancer Center Amsterdam, Cancer Biology & Immunology, Amsterdam, the Netherlands.*

16 ⁸*Current Address: Guangzhou Laboratory, Guangdong, China*

17 **Joint senior author*

18

19 Correspondence to R. Thijssen; r.thijssen@amsterdamumc.nl at the Department of Hematology,
20 Amsterdam UMC, De Boelelaan 1117, Amsterdam, The Netherlands.

21

22 Running title: scRaCH-seq to detect isoforms and mutations

23

24

25

26 **Abstract**

27

28 Single-cell long-read sequencing has transformed our understanding of isoform usage and the
29 mutation heterogeneity between cells. Despite unbiased in-depth analysis, the low sequencing
30 throughput often results in insufficient read coverage thereby limiting our ability to perform mutation
31 calling for specific genes. Here, we developed a single-cell **Rapid Capture Hybridization sequencing**
32 (scRaCH-seq) method that demonstrated high specificity and efficiency in capturing targeted
33 transcripts using long-read sequencing, allowing an in-depth analysis of mutation status and transcript
34 usage for genes of interest. The method includes creating a probe panel for transcript capture, using
35 barcoded primers for pooling and efficient sequencing via Oxford Nanopore Technologies platforms.
36 scRaCH-seq is applicable to stored and indexed single-cell cDNA which allows analysis to be
37 combined with existing short-read RNA-seq datasets. In our investigation of *BTK* and *SF3B1* genes in
38 samples from patients with chronic lymphocytic leukaemia (CLL), we detected *SF3B1* isoforms and
39 mutations with high sensitivity. Integration with short-read scRNA-seq data revealed significant gene
40 expression differences in *SF3B1*-mutated CLL cells, though it did not impact the sensitivity of the
41 anti-cancer drug venetoclax. scRaCH-seq's capability to study long-read transcripts of multiple genes
42 makes it a powerful tool for single-cell genomics.

43

44

45

46

47

48

49

50

51

52

53 **Introduction**

54

55 Single-cell sequencing technologies have revolutionized our understanding of cell state and function
56 (Tang et al. 2009; Sandberg 2014). These approaches enable the comprehensive characterization of
57 individual cells, revealing cancer biology and tumour heterogeneity (Stewart et al. 2020; Wang et al.
58 2020; Liu et al. 2022; Mustachio and Roszik 2022; Tian et al. 2022; Wang et al. 2022; Nagler and Wu
59 2023). While single-cell RNA sequencing (scRNA-seq) has been widely used for transcriptomic
60 profiling of individual cells, it has limitations in calling mutations and quantifying isoform usage due
61 to the 5' and 3' bias induced by the fragmentation step in the library preparation protocol and the
62 widespread use of short-read technology for sequencing samples. Consequently, linking the single-
63 cell transcriptome to the mutation status of cancer cells becomes a challenge. Understanding the
64 clonal and non-clonal mechanisms of selection and adaptation in response to therapeutic pressure is of
65 importance. Moreover, the influence of isoform usage on cell states further underscores the need for
66 comprehensive methodologies. In response to these challenges, alternative methods have emerged to
67 link the transcriptome profile to isoform usage, mutations, and translocations in full-length transcripts
68 at a single-cell level (Wu and Schmitz 2023).

69

70 In recent years, several high-throughput methods have been developed to enable single-cell long-read
71 sequencing. These approaches typically involve barcoding of cDNA using existing methods such as
72 10x Genomics or Drop-seq and sequencing the indexed full-length cDNA on platforms such as Pacific
73 Biosciences (PacBio) or Oxford Nanopore Technologies (ONT) (Lebrigand et al. 2020; Joglekar et al.
74 2023). While these unbiased single-cell long-read sequencing methods can detect a larger number of
75 isoforms at a single-cell level, the lower overall sequencing level per cell reduces the ability to
76 accurately quantify isoform usage and mutation calling. To overcome this challenge, other methods
77 were developed to target genes of interest, but these methods often require primer spike-in with the
78 10x Genomics protocol and primer panel design optimization (Nam et al. 2019; Griffin et al. 2023).
79 This limits the ability to target multiple genes of interest or to perform long-read sequencing on
80 already amplified single-cell indexed cDNA.

81

82 The growing interest in high throughput single-cell multi-omic methods to enhance our understanding
83 of cancer complexity led us to develop a method that can flexibly link genotypes in single cells to pre-
84 existing short-read transcriptome data generated by various 10x Genomics protocols. CITE-seq
85 combines transcriptome with cell surface protein expression and 10x Genomics multiome combines
86 transcriptome with chromatin accessibility data (Stoeckius et al. 2017). Our single-cell Rapid Capture
87 Hybridization sequencing (scRaCH-seq) method enables the capture of multiple transcripts from pre-
88 indexed and stored cDNA independent of the 10x Genomics kit used. This approach uses biotinylated
89 probes to target genes of interest and Streptavidin beads for on-target transcript enrichment. It was
90 first detailed in Thijssen *et al.* (Blood 2022), to capture the transcriptional landscape of 17 BCL2
91 family genes in patients with chronic lymphocytic leukaemia (CLL) (Thijssen et al. 2022). That study
92 investigated the resistance mechanisms observed in chronic lymphocytic leukaemia cells after
93 treatment with the BCL2 inhibitor venetoclax treatment. Using scRaCH-seq, we successfully
94 identified the *BCL2 G101V* mutation, a novel *PMAIP1* transcript and a unique *BAX* isoform across
95 different venetoclax-relapsed samples at a single-cell level. Despite the first application of scRaCH-
96 seq to that research, a detailed method description, data analysis and method evaluation has not been
97 published.

98

99 Here, we provide a comprehensive description of the scRaCH-seq method in which we target the
100 transcripts of 2 genes *SF3B1* and *BTK* as illustrative examples. The aim of this study is to demonstrate
101 the utility and accuracy of scRaCH-seq in detecting gene-specific mutations and splicing events in
102 single cells, with a focus on probe capture efficiency, artefact diagnostics and mutation/deletion
103 calling.

104

105 **Results**

106 **Assessing data quality and probe-capture efficiency**

107 A step-by-step guide to probe panel design (**Supplemental Fig S1**) was created to facilitate the
108 application of scRaCH-seq together with a tailored analysis pipeline that extends the *FLAMES* (Full-
109 Length Analysis of Mutations and Splicing) R package (Tian et al. 2021) and incorporates additional
110 read integrity checks and steps to remove background noise. After scRaCH-seq was performed
111 (**Figure 1**), a total of 43,762,097 raw reads were acquired through Nanopore sequencing across 21
112 samples, with a read length distribution spanning from 0 to 4,000bp (**Figure 2A**). The *FLAMES*
113 demultiplexing process was subsequently able to recover 28,573,271 reads that could be attributed to
114 specific cells. Despite an overall 35% read loss during this step, the length distribution of the retained
115 reads remained unaltered (**Figure 2A**). After demultiplexing, sufficient reads (~1,000,000) were
116 preserved for each sample (**Figure 2B**). We also tested Flexiplex (Cheng et al. 2024) and the reads
117 retained for downstream processing were similar. Next demultiplexed reads without TSO sequences
118 and poly(A) tails were discarded following the quality control step which eliminated around 49% of
119 reads, primarily those with lengths < 1500bp (**Figure 2A**). A total of 14,442,490 reads were retained
120 for subsequent read alignment and feature counting. Despite some read loss during demultiplexing
121 and read integrity checks, we achieved an average saturation rate exceeding 75% for the target genes
122 across all samples (**Supplemental Fig S2A-B**). Additionally, there was an increase in UMIs as more
123 transcripts were sequenced for each sample (**Supplemental Fig S2C-D**). scRaCH-seq exhibited lower
124 read loss during both demultiplexing and read integrity checks compared to single-cell full-length
125 transcript sequencing (scFLT-seq), an unbiased single-cell long-read sequencing method
126 (**Supplemental Fig S2E-G**).

127

128 scRaCH-seq and our *FLAMES* pipeline effectively captured the majority of cell barcodes identified by
129 *Cell Ranger* in scRNA-seq data (**Figure 2B**). After quality control, 92.6% of all the reads retained
130 aligned successfully to the targeted genes *SF3B1* and *BTK* (**Figure 2B, Supplemental Fig S3A**).
131 Duplicated reads were collapsed based on isoforms. In concordance with on-target reads, a high

132 number of UMIs were associated with *SF3B1* and *BTK*, with the majority of captured cells harbouring
133 the target genes (**Figure 2C**). Most of the off-target UMIs arose from single reads (**Supplemental Fig**
134 **S3B**), indicating the presence of low levels of background signal in the data. The captured off-target
135 genes include *HSPD1*, which is near *SF3B1*, and *TIMM8A*, located in close proximity to the *BTK*
136 gene (**Figure 2D**). These were likely read-through transcripts. Additionally, some other top off-target
137 transcripts, such as *RPL13*, were captured due to shared sequences with the target transcripts
138 (**Supplemental Fig S3C**). However, another group of off-target transcripts, such as *CD74*, were
139 neither located near the target genes nor shared similar sequences with them. This type of off-target
140 transcript is likely attributable to sequencing background. Per-cell UMI counts for the target genes in
141 the corresponding short-read data were evaluated, revealing a higher per-cell UMI capture in the
142 scRaCH-seq data for *SF3B1* and *BTK* (**Figure 2E**). In a separate scRaCH-seq experiment targeting 17
143 BCL2 family genes in CLL cells, similar UMI counts were observed in scRaCH-seq data compared to
144 short-read data (**Supplemental Fig S4A**). In another experiment using acute myeloid leukaemia
145 (AML) cells with similar read length distribution as the CLL samples (**Figure 2A, Supplemental Fig**
146 **S4B**) we successfully captured the transcripts of 24 genes (**Supplemental Fig S4C**). The AML
147 sample was processed using the 10x Genomics Chromium single-cell 3' kit. These findings indicate
148 the high efficiency of scRaCH-seq in capturing the transcripts of target genes observed in matched
149 single-cell 5' and 3' short-read datasets via long-read sequencing.

150

151 Next, we compared the gene expression of *SF3B1* and *BTK* in the scRaCH-seq data with the short-
152 read dataset. Similar to the short-read gene expression data, scRaCH-seq showed that *BTK* expression
153 was detected in the CLL/B cell population, whereas *SF3B1* was expressed across all cell types
154 (**Figure 2F-H**)

155

156 **Isoform detection by scRaCH-seq**

157 Given that scRaCH-seq operates by capturing transcripts of target genes with probes, it enables the
158 exploration of isoform usage for the targeted genes. In our dataset, we incorporated the CLL and B
159 cells from matched samples from CLL patients at diagnosis, after venetoclax relapse, in patients who

160 relapsed and subsequently received a BTK inhibitor (BTKi) and healthy donor samples (Thijssen et
161 al. 2022). The primary *BTK* isoform (*BTK-201*, ENST00000308731) captured by scRaCH-seq is a
162 protein-coding isoform with 19 exons, exhibiting high expression across all samples (**Figure 3A**). The
163 other four transcripts among the top 5 captured isoforms by scRaCH-seq are novel isoforms not
164 catalogued previously (**Figure 3A**). Transcripts #2 and #3 lack the last four exons compared to *BTK-*
165 *201*, with transcript #2 having a 16-base pair longer exon 10 compared to transcript #3. Novel
166 transcript #4 lacks exon 10 compared to the *BTK-201* transcripts. While all five *BTK* transcripts
167 detected by scRaCH-seq have a poly(A) tail in the 3' UTR, transcript #5, comprising only three exons,
168 has a different starting point at the 5' UTR. Consequently, it remains uncertain whether *BTK* transcript
169 #5 is an artefact resulting from truncation at the 5' end. No evidence of differential *BTK* transcript
170 usage was found across the different sample groups (**Figure 3A**). Among the *BTK* transcripts, four of
171 them shared a nearby TSS with the primary *BTK* isoform, except for *BTK* transcript #5 (**Figure 3A**).
172 However, *BTK* transcript #5 also shared a nearby TSS with another annotated isoform (**Supplemental**
173 **Fig S5A**). Additionally, TSS signals were detected in this region in the CAGE database (FUNTOM6)
174 (**Supplemental Fig S5B**).

175

176 In our investigation of the *SF3B1* isoforms, a previously unidentified *SF3B1* transcript emerged as the
177 predominant form, characterized by the absence of the last 14 exons in the 3' end (**Figure 3B**). Given
178 that this transcript is not the canonical *SF3B1* transcript expressed in the majority of cells, including
179 healthy donor samples, we undertook a comprehensive assessment to determine its authenticity. This
180 transcript does not appear to be an artefact induced by Nanopore sequencing, since it initiates at exon
181 1 and concludes with a poly(A) tail. To further validate its presence, we examined the distribution of
182 *SF3B1* reads in scRaCH-seq data and compared it to matched bulk RNA-seq and scFLT-seq data for
183 sample CLL2-RB. A consistent pattern in the *SF3B1* read distribution between scRaCH-seq and
184 scFLT-seq was observed, with a higher coverage of exons 1–11 compared to that of exons 12–25
185 (**Supplemental Fig S6A**). However, when comparing single-cell sequencing data (scRaCH-seq and
186 scFLT-seq) to bulk RNA-seq, bulk RNA-seq data exhibited a relatively stable read coverage across all
187 exons (**Supplemental Fig S6B**). This hints at a potential bias introduced by the 10x Genomics

188 protocol into the composition of *SF3B1* transcripts. Upon closer examination, we discovered that 86%
189 of single-cell reads mapped to exon 11 contained three cytosine (C) to thymine (T) mismatches
190 following the drop in read coverage, indicating that these reads were created due to the unexpected
191 binding of 10x poly(A) primers (**Supplemental Fig S6C**). This artefact is reflected in the read length
192 distribution plot with a higher proportion of reads shorter than the canonical isoforms of *SF3B1* (**Fig.**
193 **2A**). These artefacts were subsequently removed from the scRaCH-seq data, resulting in the
194 elimination of the novel *SF3B1* transcript which lacks the last 14 exons (**Supplemental Fig S6D**).

195

196 After removal of the artefacts, high expression levels of *SF3B1* transcript #1 (*SF3B1-211*,
197 ENST00000487698), transcript #2 (a novel transcript containing exons 1 – 16) and transcript #3
198 (*SF3B1-201*, ENST00000335508) were consistently observed across multiple samples (**Figure 3C**).
199 The expression of transcripts #4 (a novel transcript) and #5 (*SF3B1-204*, ENST00000414963) was
200 relatively low (**Figure 3C**). In contrast to *BTK*, a clear difference in *SF3B1* isoform usage was evident
201 between CLL cells from patient samples at screening and those from other CLL stages or healthy
202 donor B cells (**Figure 3C**). Specifically, *SF3B1* transcript #1 was highly expressed in screening
203 samples, whereas transcripts #2 and #3 were highly expressed in venetoclax-relapsed and healthy
204 donor samples. In contrast to the CLL and B cells, *SF3B1* transcript #3 was the dominant transcript
205 expressed by the T cells from the CLL patients at different stages and healthy donors (**Supplemental**
206 **Fig S7**). To investigate why the screening samples exhibit higher expression of transcript #1, we
207 linked isoform usage with gene expression. The CLL and B cells from all samples were isolated and
208 re-clustered, revealing distinct clusters of the different stages, including one B cell cluster (C9), four
209 screening clusters (C1, C3, C4, C12), three relapsed clusters (C0, C2, C8), and five clusters shared by
210 screening and relapsed CLL cells (C5, C6, C7, C10, C11) (**Figure 3D**). *SF3B1* transcript #1 was
211 specific to the screening CLL cells from cluster 3 (**Figure 3D**). Meanwhile, transcript #3, the
212 canonical *SF3B1* isoform, was expressed by all CLL and B cells (**Figure 3D**). Subsequently, single-
213 cell differential gene expression analysis between all the screening clusters (C1, C3, C4, C12) was
214 performed. For this analysis, other clusters were excluded to eliminate confounding factors related to

215 treatment stages. CLL cells in C3 exhibited a high expression level of *DDX17*, a gene involved in
216 almost all RNA metabolism processes(Xu et al. 2022) and *PARP15*, a negative regulator of
217 transcription (Ryu et al. 2015) (**Figure 3E**).

218

219 **Calling *SF3B1* mutations using scRaCH-seq data**

220 Besides isoform detection, scRaCH-seq has the ability of linking the mutation status to cell state and
221 differential gene expression. Among the CLL patients treated with venetoclax, *SF3B1* mutation was
222 detected in 5 samples by whole exome sequencing (WES) (Thijssen et al. 2022) which was used as a
223 ground truth to validate the scRaCH-seq data. For the identification of *SF3B1* mutations, we
224 aggregated the scRaCH-seq data from all the samples, revealing many *SF3B1* transcripts with single-
225 nucleotide polymorphism (SNP) (**Figure 4A**). The three most prominent SNPs were #5, #6, and #8.
226 While SNP #6 and #8 exhibited a low frequency (~20%) across all samples, SNP #5 was unique to
227 four specific samples. SNP #5 corresponds to the SF3B1 K700E mutation (Chr2:197402110 T→C), a
228 mutation frequently observed at a sub-clonal level in CLL patients (Wan and Wu 2013; Landau et al.
229 2015; Landau et al. 2017). The K700E mutation can interrupt the recruitment of *SF3B1* to the correct
230 3' end branch site and consequently result in alternative splicing (Zhang et al. 2019). We identified the
231 SF3B1 K700E mutation in CLL3 and CLL26 in both screening (S) and venetoclax-relapsed (R)
232 samples (**Supplemental Fig S8A**). This mutation was also confirmed by the WES data from these
233 samples(Thijssen et al. 2022). In scRaCH-seq, the K700E mutation was consistently observed with an
234 overall frequency of 45.0% across the four samples, encompassing 87,674 UMI counts at the location
235 of the SF3B1 K700E mutations and 48,224 UMI counts containing the reference base (Thymine)
236 (**Figure 4B**). Although SNP #1 and #3 were specific to four samples and occurred at a high frequency,
237 the overall UMI coverage for these altered transcripts was low (**Figure 4A-B**). A comparison of
238 scRaCH-seq efficiency in capturing the SF3B1 K700E mutation with whole transcriptome short-read
239 sequencing (scRNA-seq) and long-read sequencing (scFLT-seq) revealed that scRaCH-seq
240 significantly increased the capture of cells carrying the mutation and enriched for transcripts (**Figure**

241 4C). scRaCH-seq obtains a 14-fold enrichment of reads for SF3B1 K700E compared to unbiased
242 scFLT-seq, while approximately 168,000 cells were sequenced versus 1,300 cells.

243 Subsequently, we isolated cells harbouring SF3B1 K700E mutant reads from the scRaCH-seq data
244 and visualized their distribution on the UMAP layout, which contained all patient sample cells
245 (**Figure 5A**). Most cells with the mutant reads clustered within the CLL-cell group, consistent with
246 the findings from WES. However, a small number of cells with mutant reads were detected within the
247 non-CLL clusters, including the B cell cluster from a healthy donor, T cell cluster, and monocyte
248 cluster (**Figure 5A and Supplemental Fig S8B**). Given that WES established the absence of *SF3B1*
249 mutations in non-CLL cells, we investigated whether these non-CLL cells with *SF3B1* mutant reads
250 might be attributed to false positives induced by doublets. Comparative analysis of UMI counts and
251 gene counts revealed that non-CLL cells with *SF3B1* mutant reads exhibited similar values to CLL
252 cells with *SF3B1* mutants (**Figure 5B**). Moreover, expression of multiple lineage markers within the
253 non-CLL cells with *SF3B1* mutant reads was not observed, suggesting that these cells were not
254 doublets (**Figure 5B**). In the *SF3B1* wild-type samples confirmed by WES, we identified altered reads
255 with a read fraction of less than 1% (**Figure 5C**), indicative of reads containing sequencing errors
256 (Sereika et al. 2022). To mitigate these false positives in both wild-type samples and non-CLL
257 clusters, we established a threshold based on the abundance of *SF3B1* mutant transcripts detected per
258 cell. This resulted in the exclusion of 80.7% of cells with *SF3B1* mutant transcripts in wild-type
259 samples (**Figure 5C**). Furthermore, by raising the threshold to more than 2 *SF3B1* mutant transcripts,
260 93.4% of false positives in wild-type samples were successfully eliminated (**Figure 5C**). Consistent
261 with false positives observed in wild-type samples, most *SF3B1*-mutant cells in non-CLL clusters
262 exhibited low counts of altered reads (**Figure 5D**). Applying the threshold of more than 2 mutant
263 reads per cell for the identification of a true positive mutant *SF3B1* cell resulted in the removal of
264 most false positives in non-CLL clusters (**Figure 5E**). Employing the threshold of more than 2 mutant
265 transcripts based on abundance was then established as a standard quality control step to eliminate
266 false positives during mutation calling. The transcript with altered *SF3B1* #1 was observed in matched
267 samples from CLL3 and CLL5 (**Supplemental Fig S8C**). This alteration was observed in both CLL

268 clusters and non-CLL clusters in the UMAP (**Supplemental Fig S8D**). These observations suggest
269 that the C-to-G alteration in the *SF3B1* 3' UTR is likely a SNP present at a baseline level in specific
270 patient samples. Notable, even with the threshold in place, we identified SNP #8, a C→T alteration at
271 Chr2:197421097, in all samples and across all cell types (**Supplemental Fig S8E-F**).

272 **Detecting the 6-base pair deletion in *SF3B1* among CLL cells**

273 In addition to the *SF3B1* K700E point mutation, WES also confirmed the presence of a 6-bp deletion
274 (6bp-del) at Chr2:197402104:197402109, resulting in an *SF3B1* KVR700**R mutation in CLL17 at a
275 subclonal level (Thijssen et al. 2022). The original *FLAMES* pipeline was only designed to identify
276 point mutations, so to address this limitation, we developed a supplementary script dedicated to detect
277 the specific deletion directly in the FASTQ files of all samples. To minimize potential false positives,
278 we applied the same criterion as we did for point mutation calling, necessitating at least two
279 transcripts with the deletion per cell across all samples. Notably, cells with the 6bp deletion were
280 exclusively observed in sample CLL17-R (**Figure 6A**), with a predominant clustering within the CLL
281 group (**Figure 6B**), aligning with WES findings and underscoring the reliability of the deletion
282 counting method. Subsequently, single-cell differential gene expression analysis was performed to
283 identify gene expression changes between CLL cells with the *SF3B1* deletion and wild-type cells
284 across all venetoclax-relapsed (R) samples. The analysis revealed a total of 377 DEGs with a $|\log_2 \text{FC}|$
285 > 0.5 and adjusted p -value < 0.05 . Among these, 152 genes showed significant upregulation, while
286 225 genes were downregulated in CLL cells with the KVR700**R mutation (**Figure 6C**). Similarly,
287 we identified 147 differentially expressed genes (FDR < 0.05) between K700E mutant and wild-type
288 CLL cells, comprising 28 upregulated genes and 119 downregulated genes in *SF3B1* K700E CLL
289 cells (**Supplemental Fig S9A**). Given the sub-clonal nature of the *SF3B1* KVR700**R mutation in
290 CLL17, DE analysis between mutant and wild-type *SF3B1* CLL cells of CLL17 was also conducted.
291 We did not observe any DEGs (**Figure 6D**), implying that the identified DEGs between cells carrying
292 the *SF3B1* 6-bp deletion and wild-type cells are likely relapsed patient-specific rather than mutation-
293 driven.

294

295 **Altered splicing in CLL cells with *SF3B1* mutations**

296 It has been reported that the inhibition of RNA splicing can enhance the sensitivity of venetoclax in a
297 mouse model bearing acute myeloid leukaemia (AML) (Wang et al. 2023). Therefore, the question
298 remains if *SF3B1* mutation and downstream altered splicing events can contribute to venetoclax
299 resistance in patients with CLL. The current scRaCH-seq probe panel was designed for two target
300 genes to detect isoform usage and mutation calling. However, scFLT-seq from the same samples
301 offers a broader overview of all genes at the transcript level, albeit with lower read coverage per
302 gene(Thijssen et al. 2022). Interrogating the single-cell full-length transcriptomic data showed altered
303 splicing in the CLL cells with *SF3B1* mutation including altered '3 splicing of the *TPT1* gene as an
304 example (**Supplemental Fig S10A-B**). This is coherent with previous reports (Tang et al. 2020;
305 Cortés-López et al. 2023), however, we found no evidence of altered splicing transcripts specific to
306 the relapsed *SF3B1* mutated samples (**Supplemental Fig S10C**). Furthermore, no altered splicing of
307 the *BCL2* family genes was observed in *SF3B1* mutated CLL cells. This finding, in conjunction with
308 the shared enriched pathways uncovered by scRNA-seq data (**Supplemental Fig S9B-C**), suggests
309 that *SF3B1* mutations may not play a significant role in the development of acquired venetoclax
310 resistance in CLL cells.

311

312 **Discussion**

313

314 Cell identity and function could be influenced by the alternative splicing of transcripts which will
315 result in a substantial number of transcript isoforms. However, a significant portion of alternative
316 transcripts will not be detected through high-throughput sequencing-based single-cell RNA-seq
317 methods due to the short read length and the inherent bias toward 3' or 5' ends of the transcripts
318 (Joglekar et al. 2023). To address this limitation, we developed scRaCH-seq, demonstrating high
319 specificity and efficiency in capturing targeted long-read transcripts. This method provides an in-
320 depth analysis of transcript usage of genes of interest, adding another layer to the existing single-cell
321 short-read RNA-seq data. Integration of unique barcode primers specific to ONT allowed pooling and

322 efficient sequencing of multiple samples on the Nanopore platform. scRaCH-seq's capability to
323 capture multiple transcripts simultaneously enhances isoform usage detection and facilitates the
324 discovery of transcripts with mutations, providing a valuable tool across diverse research fields. Our
325 study demonstrated the efficient capture of transcripts <5kb from pre-indexed and stored cDNA,
326 depending on which genes were targeted. While this study focused on the capture of 2 genes that are
327 highly expressed, we demonstrated the efficient capture of the transcripts of 17 and 24 genes, with the
328 potential for larger gene sets without the need for primer optimization. Notably, scRaCH-seq
329 demonstrated concordance with scRNA-seq (10x Genomics) in capturing genes, suggesting its
330 reliability in transcriptomic profiling.

331

332 We optimized the *FLAMES* pipeline for scRaCH-seq data by incorporating read integrity checks to
333 eliminate 3'-end-truncated reads and setting a quality control threshold of at least two mutant reads to
334 reduce false positives. Additionally, a deletion counting script was created for identifying known
335 deletions. This allowed the detection of the SF3B1 K700E mutation and the 6-base pair deletion in
336 *SF3B1* in particular samples, both of which were previously confirmed by WES. Meanwhile, the
337 incidence of false positives stemming from sequencing errors was minimal after the filtering based on
338 the abundance of mutant reads detected in CLL cells. While single-cell DNA sequencing, such as
339 Mission Bio Tapestry, provides insight into clonality upon drug resistance (Thompson et al. 2022), it
340 lacks transcriptome profiling per cell. In contrast, scRaCH-seq offers short-read whole transcriptomic
341 data and mutation status if the mRNA of the targeted gene is expressed. By integrating the short-read
342 scRNA-seq data, we observed multiple significant DEGs between SF3B1 K700E/KVR700**R and
343 wild-type CLL cells. The subclonal nature of SF3B1 KVR700**R in CLL17 allowed us to study the
344 mutation's impact at a single-cell level. No DEGs were observed between SF3B1 KVR700**R and
345 wild-type CLL cells from sample CLL17. However, scFLT-seq demonstrated that the SF3B1
346 KVR700**R mutation increased the expression of novel isoforms in CLL17. Our single-cell data
347 revealed that the altered gene expression between *SF3B1*-mutated and wild-type CLL cells was
348 primarily driven by patient specificity. This emphasizes the importance of caution when drawing

349 conclusions based on bulk RNA-seq data comparing wild-type and mutant samples from different
350 patients.

351 Venetoclax is an effective therapy for CLL patients with the potential to induce long remissions.
352 However, we showed that multiple mechanisms resulting in the deregulation of apoptotic genes could
353 occur at a polyclonal level in venetoclax-relapsed patient samples (Thijssen et al. 2022). Now by
354 incorporating scRaCH-seq, we can add another layer of information and study the effect of an *SF3BI*
355 mutation on venetoclax sensitivity. While it was observed that inhibiting RNA-splicing could enhance
356 venetoclax sensitivity in AML (Wang et al. 2023), we demonstrated that *SF3BI* mutated venetoclax-
357 relapsed CLL cells did not express novel isoforms of the BCL2 family members that could impact
358 venetoclax sensitivity. We observed differential transcript usage of *SF3BI* between screening and
359 venetoclax-relapsed samples. *SF3BI*-211 was found to be specific to a cluster of screening samples
360 but the specific role of this *SF3BI* isoform remains unknown. The CLL cells with this *SF3BI*-211
361 isoform expressed lower levels of ribosomal genes. These findings suggest a potential connection
362 between the usage of *SF3BI*-211 transcript and the senescence stage in CLL cells, possibly through its
363 impact on gene transcription activity and protein synthesis.

364 Lastly, we observed a consistent *SF3BI* artefact across all samples in the scRaCH-seq and scFLT-seq
365 data. Comparison of read coverage between scRaCH-seq, scFLT-seq, and bulk RNA-seq data revealed
366 that this artefact resulted from the unexpected binding of the 10x reverse transcriptase primer, rather
367 than being a technical issue with the scRaCH-seq protocol. This unexpected binding of the 10x
368 primers did not impact the results of scRNA-seq since scRNA-seq operates at the gene level, counting
369 all reads aligned to *SF3BI*. However, it could pose a problem for scFLT-seq and scRaCH-seq, which
370 distinguish reads at the transcript level. Therefore, when identifying novel transcripts in scRaCH-seq
371 or scFLT-seq data, it is advisable to conduct a thorough examination of potential unexpected binding
372 sites of 10x primers within the corresponding gene.

373 In conclusion, scRaCH-seq provides an innovative strategy for studying long-read transcripts from
374 pre-indexed cDNA, holding promise for advancing gene expression studies and unravelling complex

375 biological processes. scRaCH-seq can be scaled to the number of target transcripts. We demonstrated
376 a 14-fold read enrichment in capturing a mutation compared to unbiased long-read sequencing, while
377 21 samples were processed on a single PromethIon flow cell versus 1 sample on a flow cell for
378 unbiased long-read sequencing. Its potential adaptability for PacBio sequencing further extends its
379 utility, with the primary advantage lying in linking transcript usage and mutation status at a single-cell
380 level. The approach is both cost-effective, high throughput and flexible, which is achieved by
381 leveraging a wide range of widely used 10x single-cell protocols, making scRaCH-seq applicable for
382 large-scale studies to comprehensively characterize cellular heterogeneity. It holds potential for
383 integration into single-cell spatial data, representing a powerful advancement in single-cell genomics
384 for understanding cellular heterogeneity in development, disease, and other biological processes.

385

386 **Materials and methods**

387

388 **Samples**

389 The full-length cDNA of 18 peripheral blood samples from chronic lymphocytic leukaemia (CLL)
390 patients and 3 healthy donors (HD) were used(Thijssen et al. 2022). 10x Genomics was performed
391 using the Chromium Next GEM (Gel Bead-In Emulsions) Single Cell V(D)J (v1.1, 10x Genomics,
392 cat#PN-1000165) according to the manufacturer's instructions. The indexed full-length surplus
393 cDNA(Thijssen et al. 2022) was stored and used as input for this scRaCH-seq experiment (**Figure 1**).

394

395 **Probe panel design**

396 Probes were designed to capture the indexed cDNA of target genes. Each probe was 120 base pairs
397 (bp) in length, ensuring robust coverage of the targeted regions. First, a FASTA file was generated for
398 all the isoforms of the genes of interest. The FASTA file containing all Ensembl-annotated exons from
399 all the isoforms was compiled and this was achieved by extracting the corresponding sequences from
400 the human genome GRCh38. This FASTA file served as the foundation for probe design. Next, exon
401 sequences shorter than 120 bp were merged with the preceding and subsequent exons to create

402 concatenations exceeding 120 bp, resulting in an updated FASTA file with all reads over 120 bp.
403 Using the custom FASTA file, a probe panel was generated to cover each base in the input, employing
404 <https://sg.idtdna.com/pages/tools/xgen-hyb-panel-design-tool>. For probe selection, a strategy was
405 implemented to cover every 1000 bp of an exon sequence with a single probe (**Supplemental Fig S1**).
406 Probes were selected based on their GC content (GC%) to ensure consistent efficiency during their
407 hybridization with cDNA. The average GC% of the redundant probe panel was used as a baseline.
408 The probe with GC% closest to the average GC% was chosen from the group of probes covering a
409 thousand-base region. For concatenated exons, probes fully covering the entire sequence of the
410 original short exon were selected (**Supplemental Fig S1**). Duplicated probes in the selected panel
411 were removed to finalize the probe panel. A pre-fixed R script (Team 2021) detailing the code for the
412 probe panel design process as well as a step-by-step instruction is available at the following GitHub
413 repository: https://github.com/HongkePn/RaCHseq_Probe_Design.

414

415 **Single-cell Rapid Capture Hybridization sequencing (scRaCH-seq)**

416 *Amplification of cDNA*

417 The surplus 10x indexed cDNA was amplified to increase the cDNA yield for the hybridization step.
418 To obtain (1.5-2 µg) cDNA, 2-10ng of cDNA was used as input in 5×50µl reactions. For a 50µl
419 reaction using 10× Genomics 3' or 5' cDNA, a mixture containing 10µl 5× PrimeSTAR buffer
420 (Takara, #R050B), 4µl dNTP (2.5mM; Takara, #R050B), 1µl Partial R1 –
421 CTACACGACGCTCTCCGATCT (10µM), 1µl T5' PCR Primer IIA –
422 AAGCAGTGGTATCAACGCAGAG (10µM), cDNA (5-10ng), nuclease-free water (33µl-volume
423 cDNA), 1µl PrimeSTAR GXL DNA Polymerase (Takara, #R050B) was made. The thermocycler was
424 programmed as follows: 98°C for 30sec, 9cycles of 98°C for 10sec, 65°C for 15sec and 68°C for 8min,
425 1cycle of 68°C for 10min. The concentration of the resulting product was assessed, and if it fell below
426 10 ng/µl, additional cycles were carried out as needed. The pooled amplified cDNA was cleaned up
427 with 0.6× AMPure XP or SPRIselect beads and taken up in 50µl 5× diluted Elution Buffer (EB;
428 Qiagen) in nuclease-free water (Thermo Fisher).

429

430 We optimised scRaCH-seq so that it can also be applied to 10x Genomics multiome cDNA. A mixture
431 containing 2µl cDNA (2ng), 25µl 2× KAPA HiFi master mix (Roche, #KK2602), 1µl Partial R1 –
432 CTACACGACGCTCTTCCGATCT (10µM), 1µl TSO_{p1} – TGGTATCAACGCAGAGTACATGGG
433 (10µM), 21µl EB buffer was made. The thermocycler needs to be programmed as follows: 98°C for
434 3min, 8 cycles of 98°C for 15sec, 64°C for 20sec and 72°C for 7min, and 1 cycle of 72°C for 10min.

435

436 *Hybridization*

437 The amplified cDNA sample (1.5-2µg) was dried using a DNA vacuum concentrator (speed vac)
438 together with 1µl of 1000µM IDT Blocking Oligos: Poly (A):
439 TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTT/3InvdT/, PR1: CTACACGACGCTCTTCCGATCT,
440 SO: AAGCAGTGGTATCAACGCAGAGTAC. Subsequently, the dried sample was taken up in the
441 following mixture: 8.5µl of 2× Hybridization Buffer, 2.7µl Hybridization Buffer Enhancer (xGen IDT,
442 cat#1080577) and 1.8µl nuclease-free water. The sample was incubated at 95°C for 10min to denature
443 the cDNA. To this denatured sample, 4µl of custom probe panel was added and incubated at 65°C for
444 16 hours to facilitate probe hybridization with the cDNA.

445

446 *Probe Pulldown*

447 We used the xGen IDT Lockdown Hybridization and Wash Kit. The sample was added directly to the
448 dried washed beads (100µl of Dynabeads M-270 Streptavidin) and incubated at 65°C for 45min.
449 Every 10min the sample and beads were mixed. The captured cDNA was then thoroughly washed,
450 and the bead plus sample mixture was resuspended in 50µl EB.

451

452 *Amplification of Captured cDNA Sample*

453 The captured cDNA with beads was amplified using LA Taq DNA Polymerase Hot-Start (Takara, #
454 RR042B) in 4×50µl reactions. For a 50µl reaction, a mixture containing 5µl 10x PrimeSTAR buffer,
455 4µl dNTP (2.5mM), 1µl FPS_{filA}: ACTAAAGGCCATTACGGCCTACACGACGCTCT TCCGATCT

456 (10 μ M), 1 μ l RPSfilBr: TTACAGGCCGTAATGGCCAAGCAGTGGTATCAACGCAGAGTA
457 (10 μ M), 12.5 μ l cDNA on beads, 26.1 μ l nuclease-free water, 0.4 μ l LA Taq DNA Polymerase was
458 made. The thermocycler was programmed as follows: 95°C for 2min, 12cycles of 95°C for 20sec, and
459 68°C for 10min, 1cycle of 72°C for 10min. After amplification, the cDNA was cleaned up with 0.6 \times
460 AMPure XP or SPRIselect beads and taken up in 30 μ l EB.

461

462 **Nanopore library preparation and sequencing**

463 The Oxford Nanopore Technology (ONT) SQK-PCB111.24 kit was used to index the samples. Since
464 the cDNA will have the ONT overhang, we started from the “Selecting for full-length transcripts by
465 PCR” step in the protocol with 5 μ l (1ng/ μ l) scRaCH-seq library, 0.75 μ l Unique Barcode Primer,
466 6.75 μ l Nuclease-free water and 12.5 μ l 2 \times LongAmp Hot Start Taq Master Mix. The thermocycler
467 was programmed as follows: 95°C for 30sec, 5 cycles of 95°C for 15sec, 62°C for 15sec and 65°C for
468 1min, and 1 cycle of 65°C for 6min. After the PCR, the sample was cleaned up with 0.6 \times beads. For
469 the *BTK* and *SF3B1* gene capture, 21 samples were pooled together and sequenced on PromethION
470 R9.4.1 flow cell (**Figure 1**).

471

472 The scRaCH-seq libraries can also be sequenced on a PromethION R10.4.1 flow cell and using the
473 SQK-NBD114.24 Native Barcoding Kit 24 to index the samples and pool them together. 12 μ l of
474 200fmol of the scRaCH-seq library was combined with 0.1 μ l Diluted DNA Control Sample (DCS)
475 and the other protocol components, resulting in a 15 μ l mixture. This mixture was incubated at 20°C
476 for 10min and 65°C for 10min.

477

478 **Short-read data analysis**

479 *Data processing*

480 We used *cellranger* (v5.0.0) and *bcl2fastq* (v2.19.1) to pre-process our short-read sequencing data.
481 The percentage of mitochondrial gene counts was calculated using the ``PercentageFeatureSet``
482 function in *Seurat* (v4.0.5)(Stuart et al. 2019). This calculation targeted genes starting with the regex

483 pattern “MT-”. For each sample, the `isOutlier` function from *Scater* (v1.20.0)(McCarthy et al. 2017)
484 was used to identify low-quality cells. These were defined as cells deviating by more than 3 median
485 absolute deviations (MADs) from the median in terms of unique molecular identifiers (UMIs, both
486 higher and lower), detected gene numbers (higher and lower), and mitochondrial gene expression
487 (higher). Cells identified as low-quality were then excluded from further analysis. Library size
488 normalization and $\log_2(x + 1)$ transformation were performed using the `NormalizeData` function in
489 *Seurat*. This step was followed by the identification of highly variable genes (HVGs) using the
490 `FindVariableFeatures` function in *Seurat*, adhering to default parameters. These HVGs were then
491 scaled and centred with the `ScaleData` function in *Seurat*.

492

493 *UMAP by gene expression*

494 We performed principal component analysis (PCA) by applying the `RunPCA` function in *Seurat* to
495 scaled HVGs (n=2000) with the default parameter settings. To remove the unwanted variability due to
496 batch effects, we applied *Harmony* (v0.1.0) for batch correction (Korsunsky et al. 2019). For cell
497 clustering analysis, we employed the shared nearest neighbour (SNN) method implemented in *Seurat*.
498 To construct the SNN graph, we used the `FindNeighbors` function with the first 20 Harmony
499 corrected PCA embeddings and `k.param = 20`. For visualization, we generated a Uniform Manifold
500 Approximation and Projection (UMAP) using the `RunUMAP` function with the first 20 *Harmony*-
501 corrected PCs. To identify cell clusters, we executed the `FindClusters` function on the SNN graph,
502 employing the original Louvain algorithm with default parameters. To identify the doublets, we
503 removed the cluster that exhibited the expression of multiple immune cell lineage markers and high
504 RNA contents. Finally, we applied the same *Seurat* functions on the filtered data to re-cluster cells
505 and update the UMAP.

506

507 *Single-cell differential expression analysis*

508 To perform single-cell differential expression analysis, we employed the `FindMarkers` function of
509 *Seurat*, using the parameters configured as `test.use = “MAST”`, `logfc.threshold = 0` and `min.pct = 0`.

510

511 *Pseudo-bulk differential expression analysis*

512 The count matrix for gene expression was aggregated using the ``aggregateAcross`` function
513 incorporated in *Scater*. Subsequently, the Likelihood ratio test, implemented in *edgeR* (v3.34.0) was
514 applied to identify differentially expressed genes (DEGs)(Robinson et al. 2010).

515

516 *Gene set enrichment analysis*

517 We downloaded the C2 canonical pathways and Hallmark pathway collections using *msigdbR*
518 (v7.5.1)(Liberzon et al. 2015). All genes were then ordered in a decreasing fashion based on their log-
519 fold changes. This ordered list of genes was then used as input for the ``GSEA`` function in
520 *ClusterProfiler* (v4.4.4), which was tested against the gene set collections using default
521 parameters(Wu et al. 2021).

522

523 **scRaCH-seq data analysis**

524 We developed a customized *FLAMES* (v1.3.4) pipeline to conduct the scRaCH-seq analysis (**Figure**
525 **1**). A detailed description of the original *FLAMES* pipeline can be found in the Methods paper that
526 introduces single-cell Full-Length Transcriptome sequencing (scFLT-seq) (Tian et al. 2021).

527

528 *Data pre-processing (STEP 01)*

529 We base called the raw data of nanopore sequencing using the *Guppy* software (v3.1.5) with the
530 “dna_r9.4.1_450bps_sup_prom.cfg” configuration, resulting in the generation of FASTQ files.
531 Subsequently, we used the “find_barcode” function in *FLAMES* to demultiplex the FASTQ files,
532 which was accomplished by cross-referencing the cell barcodes identified in the corresponding
533 scRNA-seq data. We allowed 2 base pairs of edit distance for the cell barcode matching.

534

535 *Read integrity check (STEP 02)*

536 We added an extra step of reads integrity check to the *FLAMES* pipeline. The demultiplexed reads
537 underwent an examination to confirm the presence of all essential components. These components
538 include a cell barcode, a unique molecular identifier (UMI), a template switch oligo (TSO) sequence,
539 and a poly(A) tail at the end of the transcript. Users can specify the TSO sequence in the
540 “find_barcode” function in *FLAMES* and filter reads that are missing the specified TSO sequence.

541

542 *Point mutation calling (STEP 03)*

543 The reads were first aligned to the human genome GRCh38 (downloaded from GENCODE) using the
544 *FLAMES* function “minimap2_align”. We then used the mutation calling function in *FLAMES* (v0.1)
545 with the default parameters to identify the point mutations, by comparing the pile-up reads against the
546 gene of interest. For *SF3B*, the gene region Chr2: 197,388,515 – 197,435,079 from the human
547 genome GRCh38 was selected. For each sample, the positions with alteration frequencies ranging
548 from 10% to 90% were identified as mutations. The allele frequency of identified mutations was then
549 counted and subsequently summarized to generate a mutation count matrix, with the rows
550 corresponding to mutations and the columns named after the cell barcodes.

551

552 *Deletion calling (STEP 03)*

553 We incorporated an additional step specifically to detect multiple base pair deletions in genes of
554 interest. We first aligned the reads which passed the integrity check to the human genome GRCh38,
555 using *minimap2* with parameters set as “-ax splice --junc-bonus 1 -k14 --secondary=no --junc-bed”(Li
556 2018). We then compared the pile-up from the reads to the reference by checking the alignments at the
557 specific chromosomal position. For the *SF3B1* 6bp deletion previously identified by WES, read
558 alignments were checked at the Chr2:197,402,104 - 197,402,109 position to identify the presence of
559 the 6bp deletion. The reads with the multiple bp deletion were merged by UMIs and subsequently
560 counted for each cell to generate a count matrix, with rows representing the *SF3B1* 6bp-deletion/WT
561 and column names corresponding to the cell barcodes. We updated *FLAMES* to provide functions for
562 detection of multiple base pair deletions along with point mutations at both piled-up bulk level
563 (“find_variants”) and single-cell level (“sc_mutations”) to simplify such analyses.

564

565 *Isoform detection (STEP 04)*

566 We employed the *FLAMES* function “find_isoform” to summarize the alignment for making isoform
567 assembly. The isoforms that exhibited <10bp variance in their splicing junctions and <100bp variance
568 in their start or end sites were merged. Following this, we re-aligned the reads that passed the integrity
569 check to the isoform assembly to generate the isoform count matrix, with the rows corresponding to
570 isoforms and the columns named after the cell barcodes.

571

572 *Data integration (STEP 05)*

573 We combined the isoform and mutation count matrices with the processed single-cell short-read data
574 to construct a comprehensive multi-assay *Seurat* object, using the 'CreateAssayObject' function. The
575 scRaCH-seq data were used to identify cell groups with specific isoforms or mutations, followed by
576 short-read differential expression analysis between these groups.

577

578 *Differential transcript usage analysis*

579 For scRaCH-seq data, we merged the per-cell counts of target gene transcripts by samples to make
580 pseudo-bulk transcript data. We calculated the frequency of different transcripts for our genes of
581 interest *SF3B1* and *BTK* respectively and then displayed the top 5 transcripts with the highest
582 frequencies across all samples for both target genes. For scFLT-seq data, cells from patients
583 harbouring *SF3B1* mutations/deletions were combined per sample to make pseudo-bulk transcript
584 data for the *SF3B1*-mutated group. Cells of wild-type screening/venetoclax-relapsed/healthy donor
585 samples were merged to construct screening-wt/relapsed-wt/healthy-wt pseudo-bulk transcript data
586 for the *SF3B1*-wild type group. We then applied the function “diffSpliceDGE” in *edgeR* to identify
587 the differential transcript usage between groups (*SF3B1*-mutated, *SF3B1*-wt) using pseudo-bulk data
588 as replicates. The scFLT-seq data and whole exome sequencing data involved in this study can be
589 access using the EGA accession number EGAS00001005815 (Thijssen et al. 2022).

590

591 **Data Access**

592 All raw sequencing data generated in this study have been submitted to the European Genome-
593 phenome Archive (EGA; <https://ega-archive.org/>) under accession number EGAD50000000235.

594

595 The probe panel design code is available at https://github.com/HongkePn/RaCHseq_Probe_Design.

596 The complete scRaCH-seq analysis code is available at

597 https://github.com/HongkePn/scRaCHseq_data_analysis. We have also included a copy of the code in

598 the Supplemental Material of this paper (Supplemental Code).

599

600 **Competing Interest Statement**

601 The authors declare no competing interests.

602

603 **Acknowledgements**

604 We thank Stephen Wilcox, Sarah MacRaild, the WEHI SCORE team and core facilities (Genomics)

605 for their support. Illustrations were created with BioRender.com. This work was supported by

606 fellowships and grants from the Australian National Health and Medical Research Council: Program

607 Grants 1113133 to D.C.S.H.; Synergy 2011139 to A.H.W., A.W.R. and D.C.S.H.; Fellowship

608 1156024 to D.C.S.H.; Ideas Grant 2013478 to D.C.S.H. and R.T; Investigator 2017257 to M.E.R. and

609 1174902 to A.W.R, the Leukemia & Lymphoma Society of America (Specialized Center of Research

610 (SCOR) grant 7015-18 to A.W.R. and D.C.S.H.), The Australian Research Council (Discovery

611 Project 200102903 to M.E.R.), Leukaemia Foundation (grant 2012526 to R.T.), Victorian Cancer

612 Agency (ECRF21014 fellowship to R.T.), the University of Melbourne (MIRS and MIFRS

613 scholarships to H.P.), the Chan Zuckerberg Initiative DAF (an advised fund of Silicon Valley

614 Community Foundation; grant number 2019-002443 to M.E.R.), and the Amsterdam UMC

615 Fellowship (R.T.). This work was made possible through Victorian State Government Operational

616 Infrastructure Support and the Australian Government NHMRC IRIISS.

617 *Author contribution:* H.P., D.C.S.H., M.E.R. and R.T. designed the study. A.H.W. and A.W.R., were
618 responsible for patient care and recruited patient samples. H.P., J.S.J., C.C.C., N.S.A. and R.T.
619 acquired the data. H.P., L.T., C.W., Y.Y., N.A., N.M.D. and M.E.R. analysed and interpreted the data.
620 C.W. updated *FLAMES*. H.P., M.E.R. and R.T. wrote the first version of the manuscript; all authors
621 reviewed the data and contributed to critical revision of the manuscript
622

623 **References**

624

- 625 Cheng O, Ling MH, Wang C, Wu S, Ritchie ME, Goke J, Amin N, Davidson NM. 2024. Flexiplex: a
626 versatile demultiplexer and search tool for omics data. *Bioinformatics* **40**.
- 627 Cortés-López M, Chamely P, Hawkins AG, Stanley RF, Swett AD, Ganesan S, Mouhieddine TH, Dai
628 X, Kluegel L, Chen C et al. 2023. Single-cell multi-omics defines the cell-type-specific
629 impact of splicing aberrations in human hematopoietic clonal outgrowths. *Cell Stem Cell* **30**:
630 1262-1281.e1268.
- 631 Griffin GK, Booth CAG, Togami K, Chung SS, Ssozi D, Verga JA, Bouyssou JM, Lee YS,
632 Shanmugam V, Hornick JL et al. 2023. Ultraviolet radiation shapes dendritic cell leukaemia
633 transformation in the skin. *Nature* **618**: 834-841.
- 634 Joglekar A, Foord C, Jarroux J, Pollard S, Tilgner HU. 2023. From words to complete phrases: insight
635 into single-cell isoforms using short and long reads. *Transcription*
636 doi:10.1080/21541264.2023.2213514: 1-13.
- 637 Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, Baglaenko Y, Brenner M, Loh PR,
638 Raychaudhuri S. 2019. Fast, sensitive and accurate integration of single-cell data with
639 Harmony. *Nat Methods* **16**: 1289-1296.
- 640 Landau DA, Sun C, Rosebrock D, Herman SEM, Fein J, Sivina M, Underbayev C, Liu D,
641 Hoellenriegel J, Ravichandran S et al. 2017. The evolutionary landscape of chronic
642 lymphocytic leukemia treated with ibrutinib targeted therapy. *Nat Commun* **8**: 2185.
- 643 Landau DA, Tausch E, Taylor-Weiner AN, Stewart C, Reiter JG, Bahlo J, Kluth S, Bozic I, Lawrence
644 M, Bottcher S et al. 2015. Mutations driving CLL and their evolution in progression and
645 relapse. *Nature* **526**: 525-530.
- 646 Lebrigand K, Magnone V, Barbry P, Waldmann R. 2020. High throughput error corrected Nanopore
647 single cell transcriptome sequencing. *Nature Communications* **11**: 4025.
- 648 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094-3100.
- 649 Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. 2015. The Molecular
650 Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* **1**: 417-425.
- 651 Liu T, Liu C, Yan M, Zhang L, Zhang J, Xiao M, Li Z, Wei X, Zhang H. 2022. Single cell profiling of
652 primary and paired metastatic lymph node tumors in breast cancer patients. *Nat Commun* **13**:
653 6823.
- 654 McCarthy DJ, Campbell KR, Lun AT, Wills QF. 2017. Scater: pre-processing, quality control,
655 normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics* **33**: 1179-
656 1186.
- 657 Mustachio LM, Roszik J. 2022. Single-Cell Sequencing: Current Applications in Precision Onco-
658 Genomics and Cancer Therapeutics. *Cancers (Basel)* **14**.
- 659 Nagler A, Wu CJ. 2023. The end of the beginning: application of single-cell sequencing to chronic
660 lymphocytic leukemia. *Blood* **141**: 369-379.
- 661 Nam AS, Kim KT, Chaligne R, Izzo F, Ang C, Taylor J, Myers RM, Abu-Zeinah G, Brand R, Omans
662 ND et al. 2019. Somatic mutations and cell identity linked by Genotyping of Transcriptomes.
663 *Nature* **571**: 355-360.
- 664 R Core Team. 2021. R: A language and environment for statistical computing. *R Foundation for*
665 *Statistical Computing, Vienna, Austria*. URL <https://www.R-project.org/>
- 666 Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential
667 expression analysis of digital gene expression data. *Bioinformatics* **26**: 139-140.
- 668 Ryu KW, Kim D-S, Kraus WL. 2015. New Facets in the Regulation of Gene Expression by ADP-
669 Ribosylation and Poly(ADP-ribose) Polymerases. *Chemical Reviews* **115**: 2453-2481.
- 670 Sandberg R. 2014. Entering the era of single-cell transcriptomics in biology and medicine. *Nat*
671 *Methods* **11**: 22-24.
- 672 Sereika M, Kirkegaard RH, Karst SM, Michaelsen TY, Sørensen EA, Wollenberg RD, Albertsen M.
673 2022. Oxford Nanopore R10.4 long-read sequencing enables the generation of near-finished
674 bacterial genomes from pure cultures and metagenomes without short-read or reference
675 polishing. *Nat Methods* **19**: 823-826.

- 676 Stewart CA, Gay CM, Xi Y, Sivajothi S, Sivakamasundari V, Fujimoto J, Bolisetty M, Hartsfield PM,
677 Balasubramaniyan V, Chalishazar MD et al. 2020. Single-cell analyses reveal increased
678 intratumoral heterogeneity after the onset of therapy resistance in small-cell lung cancer. *Nat*
679 *Cancer* **1**: 423-436.
- 680 Stoeckius M, Hafemeister C, Stephenson W, Houck-Loomis B, Chattopadhyay PK, Swerdlow H,
681 Satija R, Smibert P. 2017. Simultaneous epitope and transcriptome measurement in single
682 cells. *Nature Methods* **14**: 865-868.
- 683 Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, Hao Y, Stoeckius M,
684 Smibert P, Satija R. 2019. Comprehensive Integration of Single-Cell Data. *Cell* **177**: 1888-
685 1902.e1821.
- 686 Tang AD, Soulette CM, van Baren MJ, Hart K, Hrabeta-Robinson E, Wu CJ, Brooks AN. 2020. Full-
687 length transcript characterization of SF3B1 mutation in chronic lymphocytic leukemia reveals
688 downregulation of retained introns. *Nat Commun* **11**: 1438.
- 689 Tang F, Barbacioru C, Wang Y, Nordman E, Lee C, Xu N, Wang X, Bodeau J, Tuch BB, Siddiqui A et
690 al. 2009. mRNA-Seq whole-transcriptome analysis of a single cell. *Nat Methods* **6**: 377-382.
- 691 Thijssen R, Tian L, Anderson MA, Flensburg C, Jarratt A, Garnham AL, Jabbari JS, Peng H, Lew TE,
692 Teh CE et al. 2022. Single-cell multiomics reveal the scale of multilayered adaptations
693 enabling CLL relapse during venetoclax therapy. *Blood* **140**: 2127-2141.
- 694 Thompson ER, Nguyen T, Kankanige Y, Markham JF, Anderson MA, Handunnetti SM, Thijssen R,
695 Yeh PS, Tam CS, Seymour JF et al. 2022. Single-cell sequencing demonstrates complex
696 resistance landscape in CLL and MCL treated with BTK and BCL2 inhibitors. *Blood Adv* **6**:
697 503-508.
- 698 Tian L, Jabbari JS, Thijssen R, Gouil Q, Amarasinghe SL, Voogd O, Kariyawasam H, Du MRM,
699 Schuster J, Wang C et al. 2021. Comprehensive characterization of single-cell full-length
700 isoforms in human and mouse with long-read sequencing. *Genome Biol* **22**: 310.
- 701 Tian Y, Li Q, Yang Z, Zhang S, Xu J, Wang Z, Bai H, Duan J, Zheng B, Li W et al. 2022. Single-cell
702 transcriptomic profiling reveals the tumor heterogeneity of small-cell lung cancer. *Signal*
703 *Transduct Target Ther* **7**: 346.
- 704 Wan Y, Wu CJ. 2013. SF3B1 mutations in chronic lymphocytic leukemia. *Blood* **121**: 4627-4634.
- 705 Wang E, Pineda JMB, Kim WJ, Chen S, Bourcier J, Stahl M, Hogg SJ, Bewersdorf JP, Han C, Singer
706 ME et al. 2023. Modulation of RNA splicing enhances response to BCL2 inhibition in
707 leukemia. *Cancer Cell* **41**: 164-180.e168.
- 708 Wang L, Mo S, Li X, He Y, Yang J. 2020. Single-cell RNA-seq reveals the immune escape and drug
709 resistance mechanisms of mantle cell lymphoma. *Cancer Biol Med* **17**: 726-739.
- 710 Wang X, Nissen M, Gracias D, Kusakabe M, Simkin G, Jiang A, Duns G, Sarkozy C, Hilton L,
711 Chavez EA et al. 2022. Single-cell profiling reveals a memory B cell-like subtype of follicular
712 lymphoma with increased transformation risk. *Nat Commun* **13**: 6772.
- 713 Wu S, Schmitz U. 2023. Single-cell and long-read sequencing to enhance modelling of splicing and
714 cell-fate determination. *Comput Struct Biotechnol J* **21**: 2373-2380.
- 715 Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, Feng T, Zhou L, Tang W, Zhan L et al. 2021.
716 clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. *Innovation (Camb)*
717 **2**: 100141.
- 718 Xu K, Sun S, Yan M, Cui J, Yang Y, Li W, Huang X, Dou L, Chen B, Tang W et al. 2022. DDX5 and
719 DDX17—multifaceted proteins in the regulation of tumorigenesis and tumor progression.
720 *Frontiers in Oncology* **12**.
- 721 Zhang J, Ali AM, Lieu YK, Liu Z, Gao J, Rabadan R, Raza A, Mukherjee S, Manley JL. 2019.
722 Disease-Causing Mutations in SF3B1 Alter Splicing by Disrupting Interaction with SUGP1.
723 *Mol Cell* **76**: 82-95.e87.

725

726

727 **Figure legends**

728

729 **Figure 1 Schematic of scRaCH-seq approach and *FLAMES* pipeline for single-cell long-read**
730 **data analysis.**

731 scRaCH-seq can be incorporated in the standard high-throughput single-cell RNA-seq experiments
732 (left side). After single-cell isolation, mRNA is converted into cDNA and barcoded. Only some of this
733 amplified indexed cDNA is used for short-read library preparation and Illumina sequencing. The
734 surplus cDNA will be stored and can be used for scRaCH-seq. For scRaCH-seq, a probe panel will be
735 designed for genes of interest. The biotinylated probes will be hybridized with amplified cDNA
736 overnight. The probes and target genes will be captured with Streptavidin beads, washed, and
737 amplified. The enriched target long-read transcripts will be sequenced on the Nanopore platform.
738 With the *FLAMES* pipeline, FASTQ files were demultiplexed by cross-referencing the cell barcodes
739 identified in scRNA-seq data (STEP 01). Next, the demultiplexed reads will undergo an integrity
740 check and the reads that possess a UMI, a TSO sequence, and poly(A) tails are retained (STEP 02).
741 Reads were aligned to the GRCh38 reference genome to construct a transcript assembly.
742 Concurrently, the piled-up reads were compared against GRCh38 to identify base alterations and
743 deletions, generating a comprehensive mutation/deletion matrix (STEP 03). The reads will be
744 realigned to the transcript assembly for quantification (STEP 04). Single-cell short-read gene
745 expression data was used to cluster the cells, based on the conventional analysis pipeline for scRNA-
746 seq. The transcript usage and mutation/deletion information were then aligned with the single-cell
747 gene expression data. The bridge connecting these datasets was the shared cell barcodes, ensuring the
748 gene expression profile at a transcript level (STEP 05).

749

750 **Figure 2 scRaCH-seq is efficient in capturing enriched genes of interest.**

751 **A)** A graph showing the losses in read counts due to demultiplexing and integrity checks, depicting
752 the distribution of read lengths ranging from 0 to 5,000 base pairs. The canonical isoform of *BTK*
753 (2,027 bp) and *SF3B1* (2,225 bp) are marked on the plot. The four shades of grey, ranging from light

754 to dark, correspond to raw reads, demultiplexed reads, reads with TSO, and reads with also a poly(A)
755 end.

756 **B)** Bar plots showing the loss of reads after demultiplexing (top panel), barcodes detected in long-read
757 sequencing data (middle panel) and counts of off-target and on-target reads for each sample (bottom
758 panel).

759 **C)** Violin plot showing the \log_{10} UMI counts of collapsed *BTK* (red), *SF3BI* (orange) and off-target
760 transcripts (purple) based on isoform detected by scRaCH-seq (left panel). Violin plot showing the
761 \log_{10} counts of cells possessing collapsed transcripts of *BTK* (red), *SF3BI* (orange) and off-target
762 genes (purple) (right panel).

763 **D)** Dot plot showing the \log_{10} counts of off-target genes with the top 10 off-target transcripts
764 specifically marked.

765 **E)** Dot plot showing the per-cell UMIs of *SF3BI* (red) and *BTK* (orange), detected by scRNA-seq (X-
766 axis) and scRaCH-seq (Y-axis) for all samples.

767 **F)** UMAP projection of peripheral blood mononuclear cells from CLL patients or healthy donors and
768 clustering based on short-read gene expression.

769 **G)** Violin plot showing the expression of CD14 (monocyte marker), CD19 (B/CLL cell marker) and
770 CD3 (T cell marker) per cell cluster (Figure 2F).

771 **H)** UMAP projection of *SF3BI* (left) and *BTK* (right) gene-level expression detected by scRNA-seq
772 (purple) or scRaCH-seq (green).

773

774 **Figure 3 scRaCH-seq reveals different *SF3BI* isoform usage by CLL cells from patients**
775 **undergoing different treatments.**

776 **A)** Heatmap showing *BTK* isoform usage (rows) per sample (columns) in the B/CLL cells. The top 5
777 *BTK* transcripts are illustrated on the right with novel transcripts highlighted in yellow. Samples are
778 grouped in screening (S; green), venetoclax relapsed (VEN-R/R; pink), venetoclax relapsed and
779 subsequently on BTKi (VEN-RB/RB; bright yellow) and healthy donor (HD; turquoise).

780 **B)** Illustration of the dominant *SF3BI* transcript identified and characterized by the absence of the last
781 14 exons in the 3' end (highlighted in purple).

782 C) Heatmap showing *SF3B1* isoform usage (rows) per sample (columns) in the B/CLL cells. The top
783 5 *SF3B1* transcripts are illustrated on the right with novel transcripts highlighted in yellow.

784 D) UMAP projection of CLL/B cells from CLL patients or healthy donors and clustering based on
785 short-read gene expression. The samples from CLL patients at screening (S; green), venetoclax
786 relapsed (R; pink), venetoclax relapsed and subsequently on BTKi (RB; bright yellow) and healthy
787 donors (HD; turquoise) are highlighted in the UMAP. The expression of the *SF3B1* transcript #1 and
788 #3 are overlaid as red dots in the UMAP (Figure 3B).

789 E) Dot plot showing marker genes that distinguish the CLL cells present in the 4 clusters shared by
790 the screening samples (C1, C3, C4 and C12 in Figure 3C). Cluster 3 with *SF3B1* transcript #1 usage is
791 indicated in bold.

792

793 **Figure 4 Mutation calling with scRaCH-seq and FLAMES.**

794 A) Dot plot showing the *SF3B1* point mutations identified using *FLAMES* (left panel). After
795 consolidating reads with UMIs, the reads from all the samples were aligned to the GRCh38 reference
796 genome. The size of each dot in the plot corresponds to the proportion of altered transcripts.

797 B) A graphical representation highlights the precise locations of these high-frequency (>10%)
798 alterations within the *SF3B1* gene. The bottom graph shows the frequency of altered *SF3B1*
799 transcripts in each sample.

800 C) Bar plots showing the cells with a *SF3B1* K700E mutation (left panel) or transcripts (UMIs) with
801 *SF3B1* K700E mutation (right panel) detected by short-read scRNA-seq, long-read scFLT-seq or
802 long-read scRaCH-seq.

803

804 **Figure 5 The *SF3B1* K700E mutation is detected by scRaCH-seq.**

805 A) UMAP projection of cells carrying the *SF3B1* K700E mutation (red).

806 B) Violin plot showing the UMI counts, gene counts (Features), expression of mitochondria genes
807 (MT%), expression of CD3 (T cell marker), CD14 (monocyte marker) and CD19 (B cell marker) per
808 cell type cluster.

809 C) Bar plot showing the abundance of mutant transcripts. Left plot showing cells per sample with >0
810 transcripts (UMI) of mutant SF3B1 K700E. Middle plot showing cells per sample with >1 transcripts
811 (UMIs) of *SF3B1* mutation. Right plot showing cells per sample with >2 transcripts (UMIs) of *SF3B1*
812 mutation. Samples with confirmed SF3B1 K700E mutation by WES are highlighted in dark grey.

813 D) Bar plot showing the distribution of the abundance of SF3B1 K700E mutation detected in CLL
814 (pink) and non-CLL (blue) cells. The vertical dashed line indicates the criterion of >2 *SF3B1* mutation
815 transcripts.

816 E) UMAP projection of cells carrying >2 SF3B1 K700E mutation transcripts (red).

817

818 **Figure 6 The SF3B1 KVR700**R alteration is detected by scRaCH-seq.**

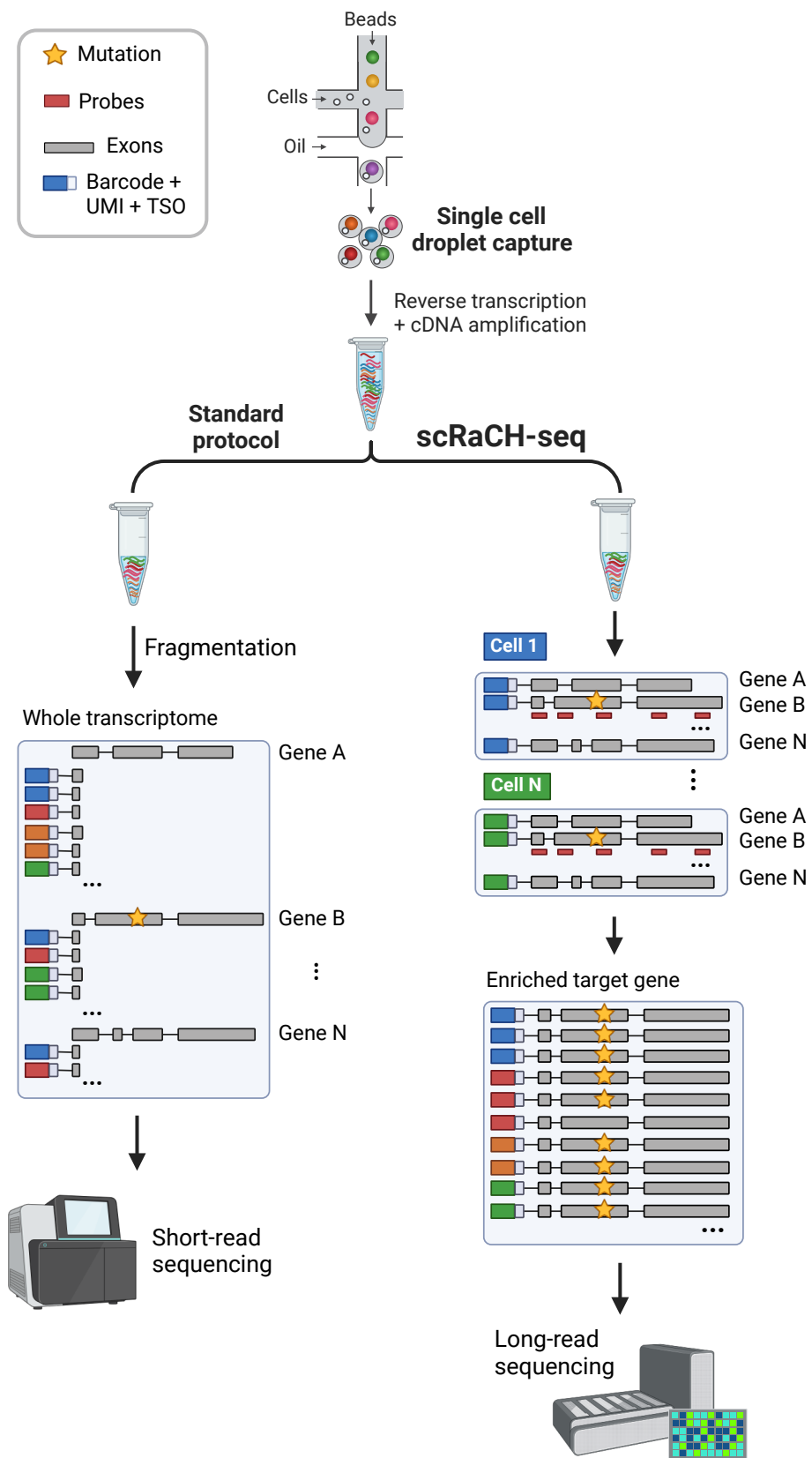
819 A) Bar plot showing the cells with *SF3B1* 6bp-deletion (Chr2:197402104:197402109) across all
820 samples. A criterion of > 2 of *SF3B1* with 6bp-del transcripts (UMIs) per cell was employed.

821 B) UMAP projection of cells carrying >2 SF3B1 KVR700**R altered transcripts (red).

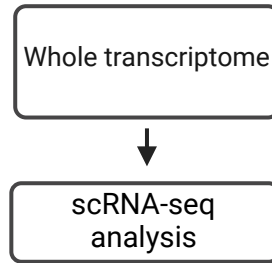
822 C) Volcano plot showing the differentially expressed genes (DEGs; FDR<0.05) between SF3B1
823 KVR700**R mutant and wild-type CLL cells from all venetoclax relapsed CLL samples.

824 D) Volcano plot showing the differentially expressed genes (DEGs; FDR<0.05) between SF3B1
825 KVR700**R mutant and wild-type CLL cells from CLL17.

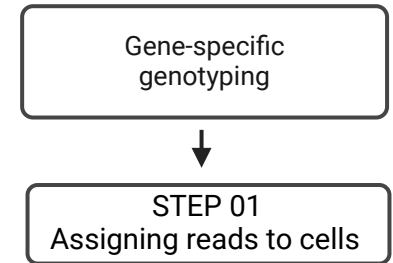
826



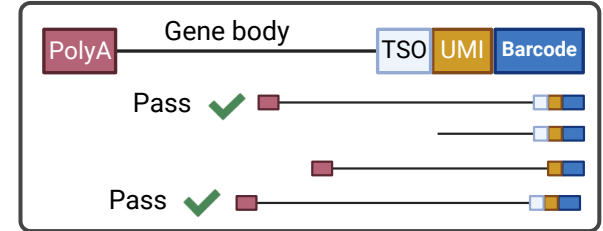
Short-read data



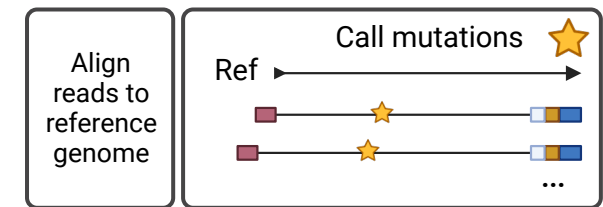
Long-read data



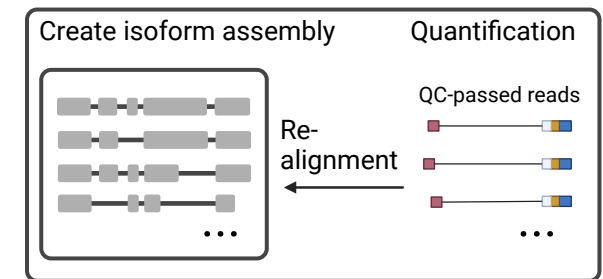
STEP 02 QC Check reads' integrity



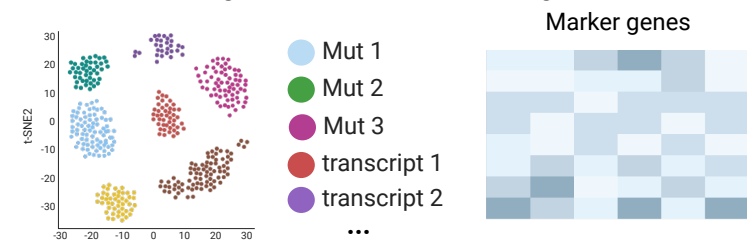
STEP 03 Mutation calling

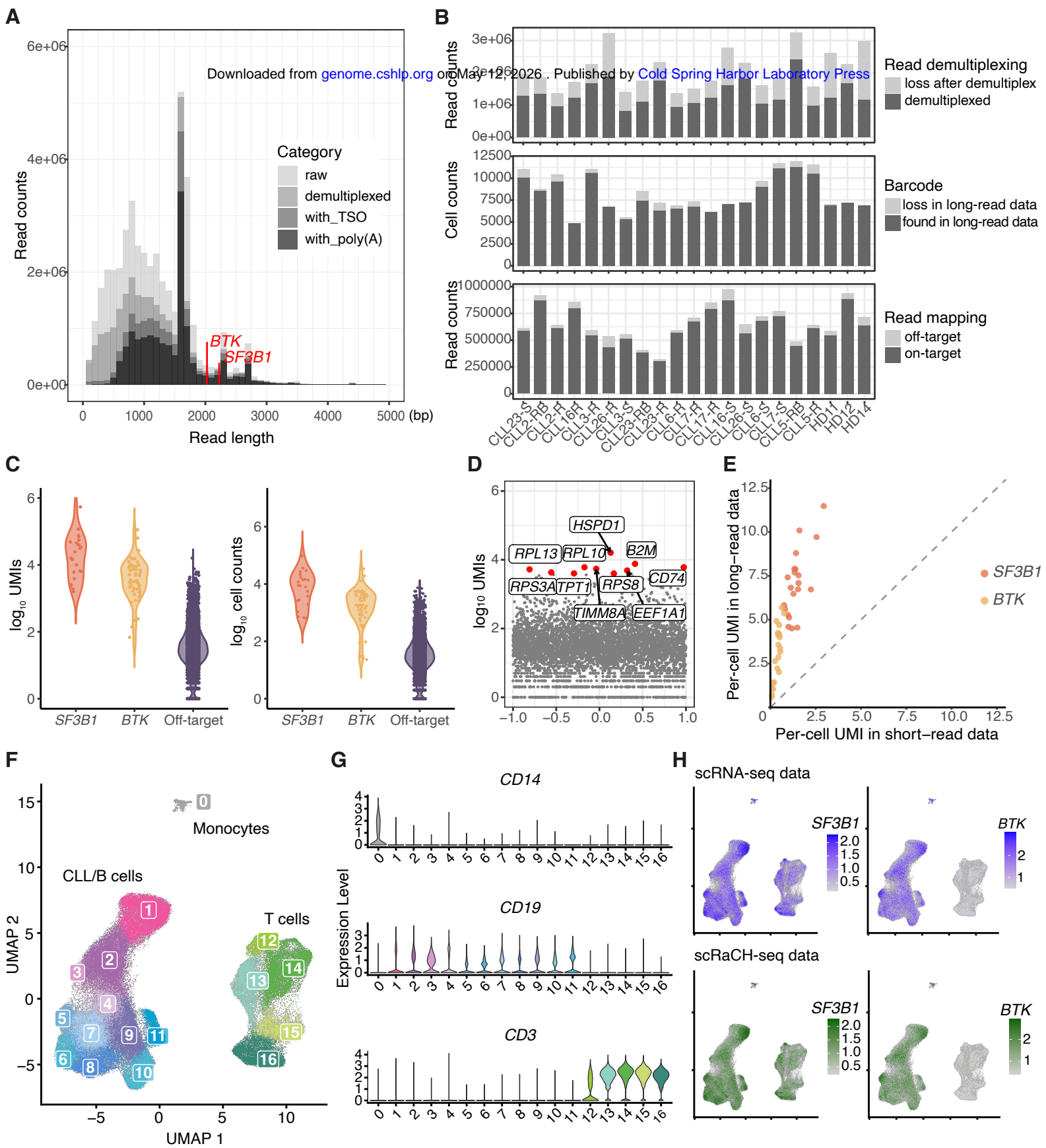


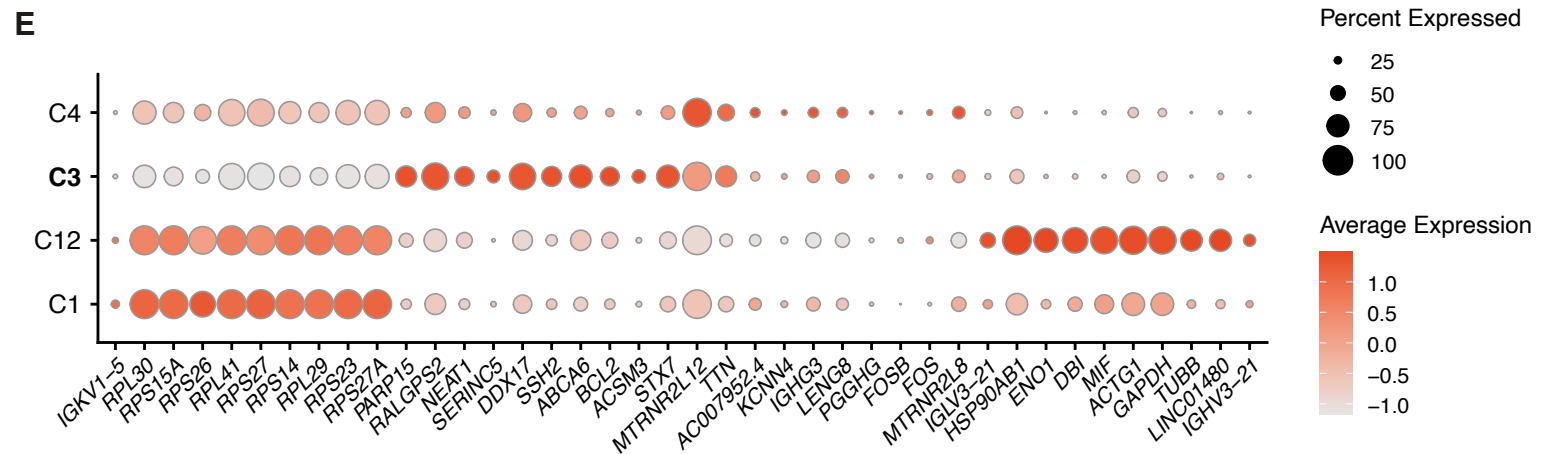
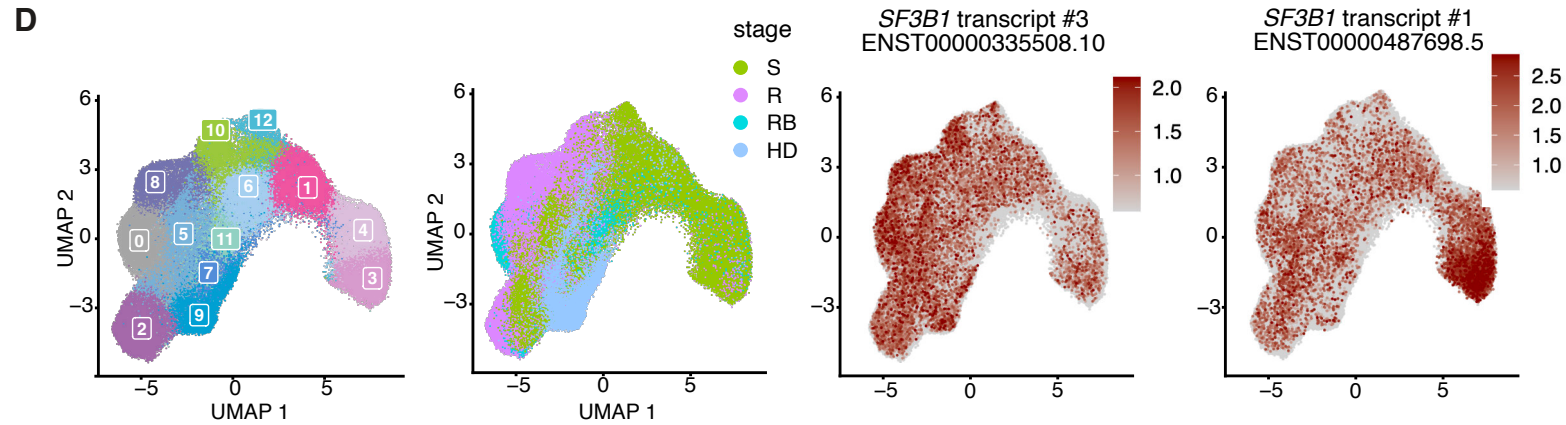
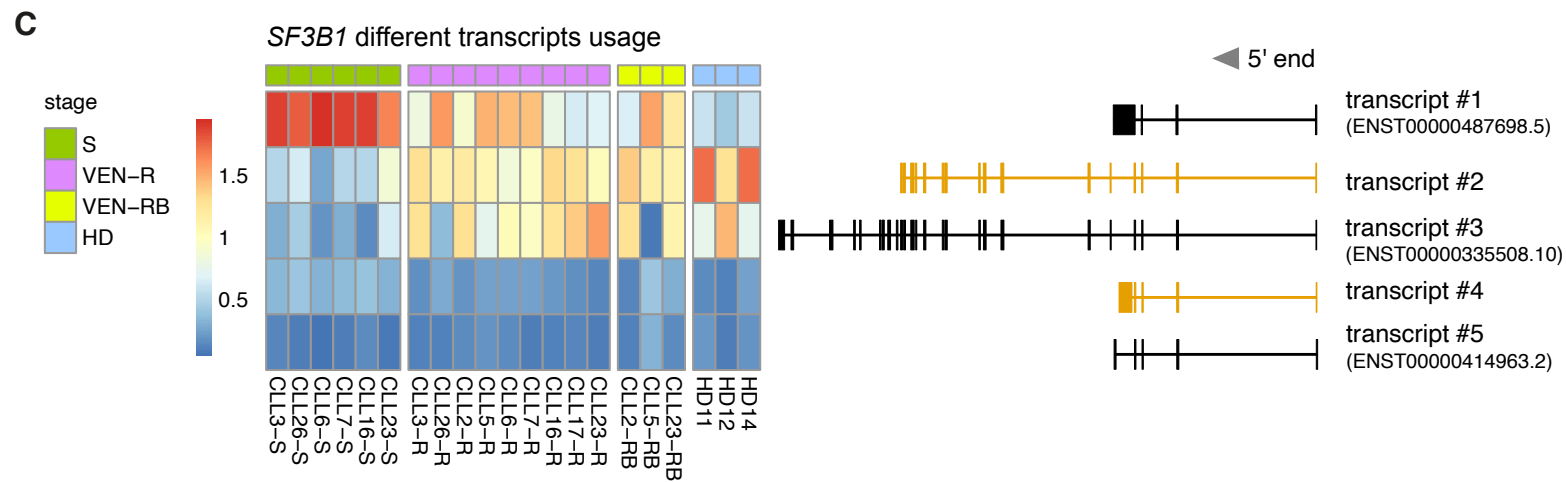
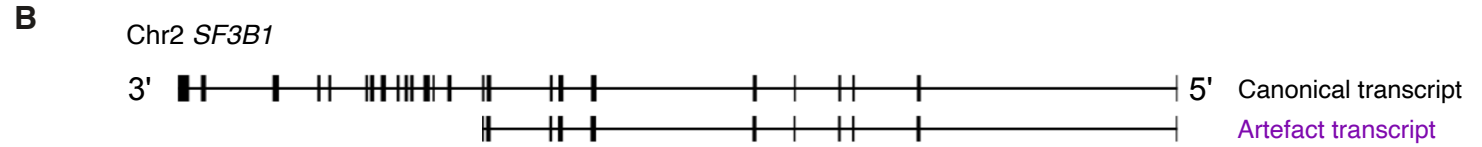
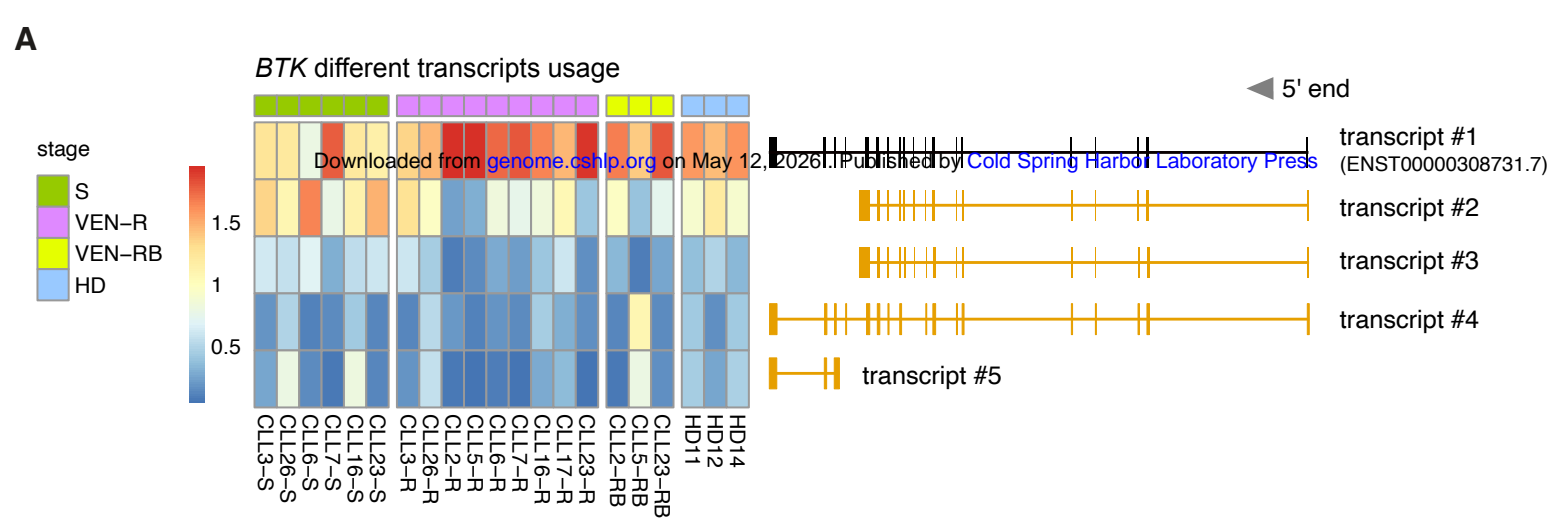
STEP 04 Isoform detection

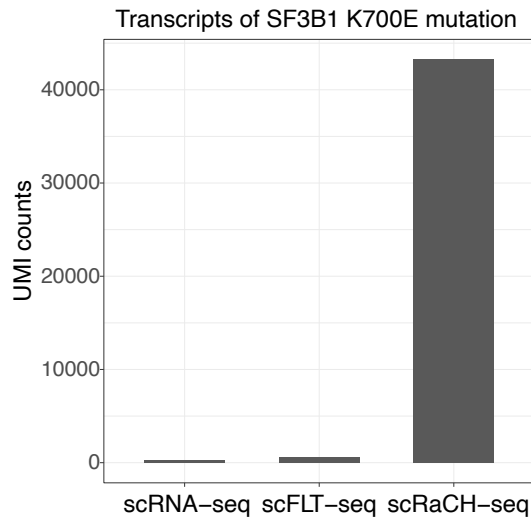
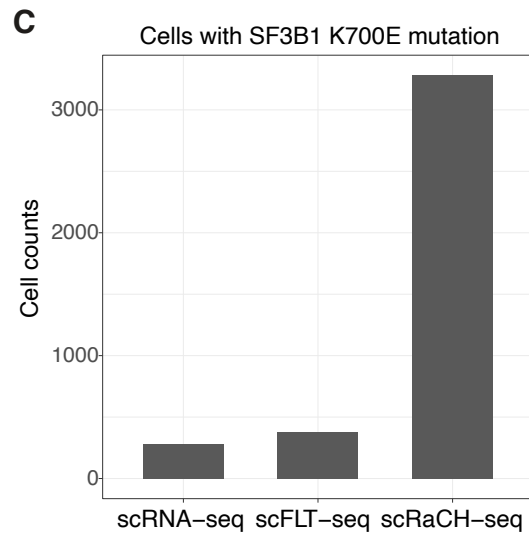
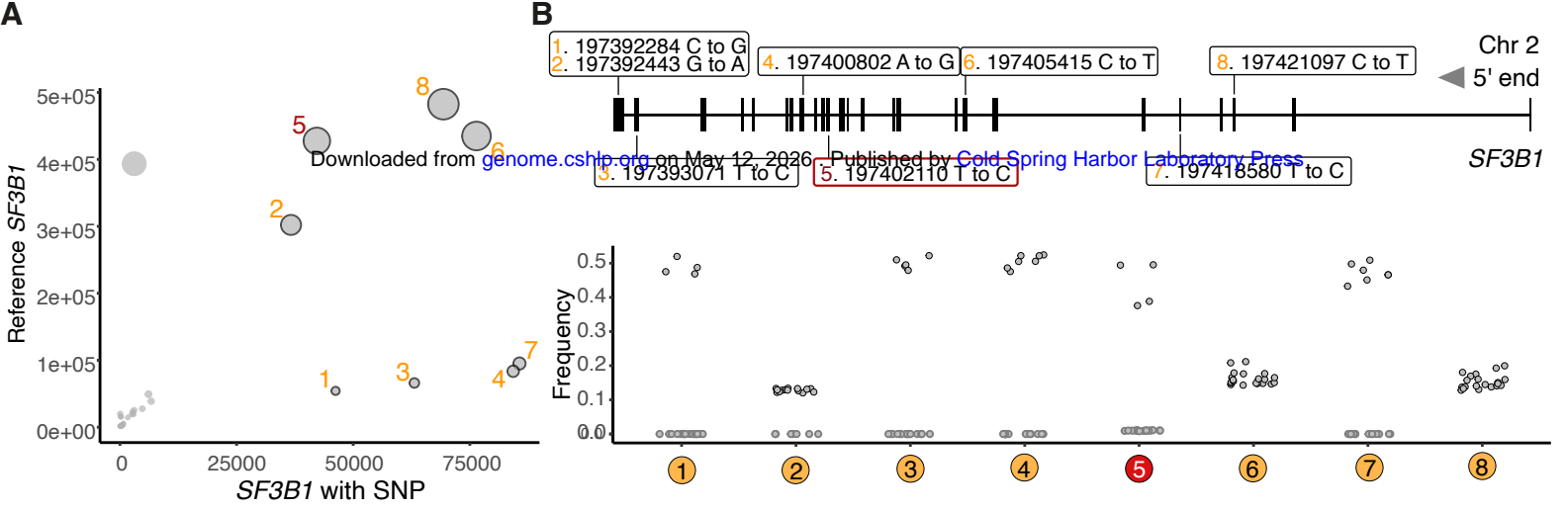


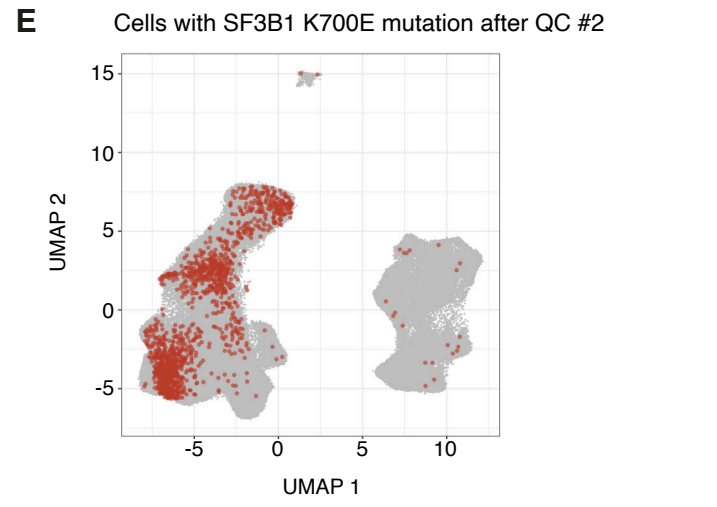
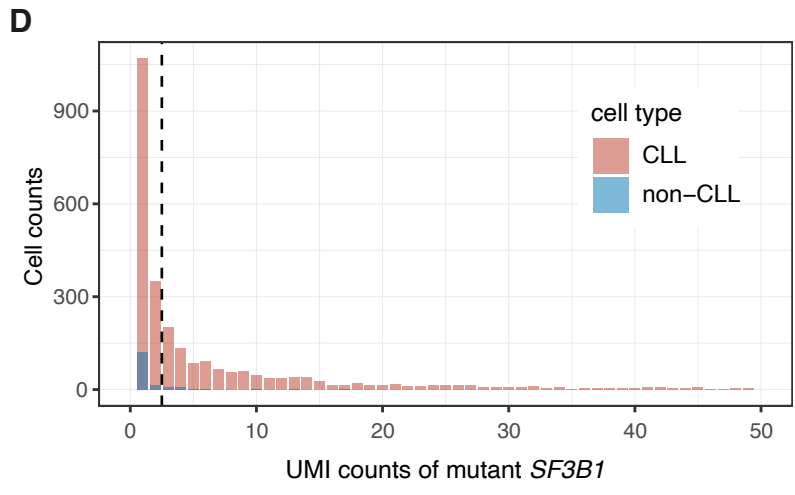
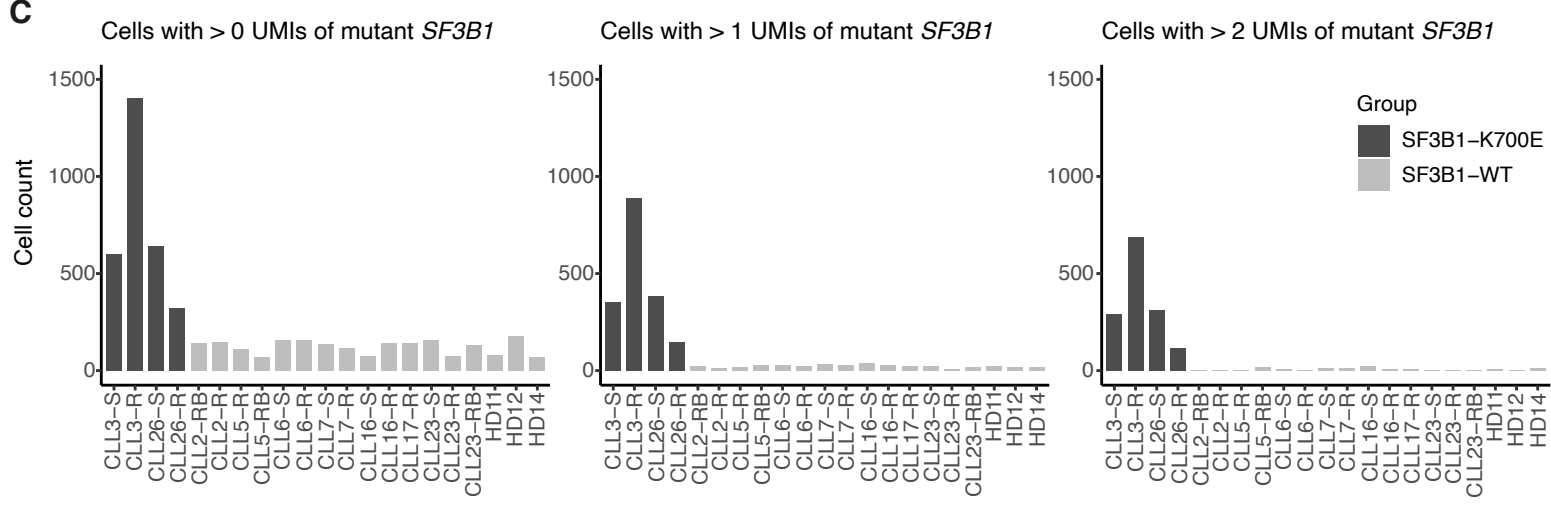
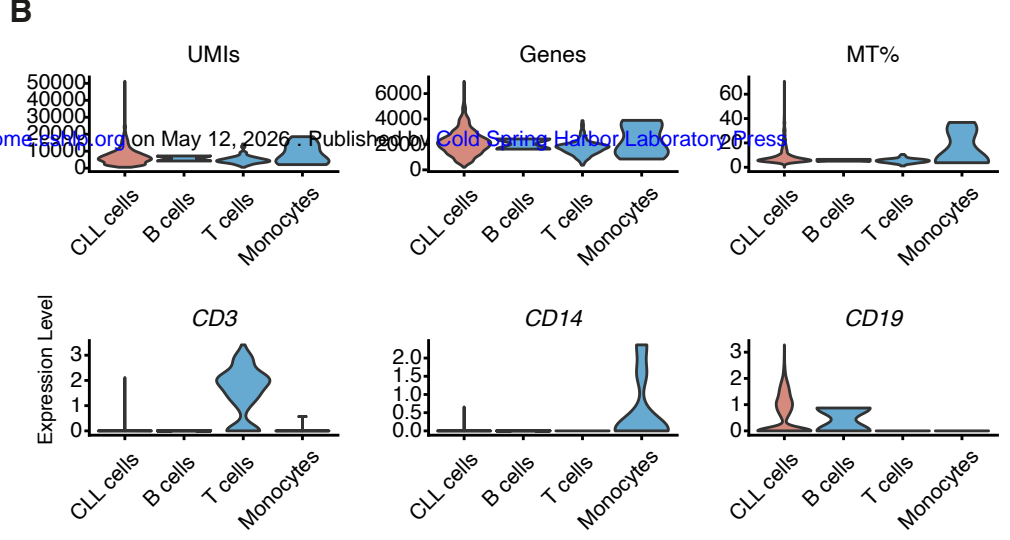
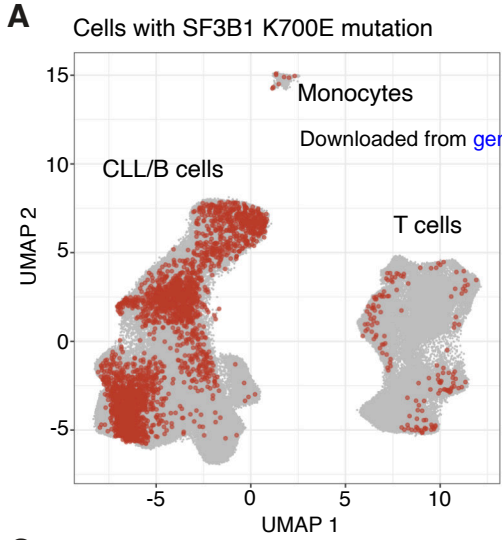
STEP 05 Long-/short-read data integration



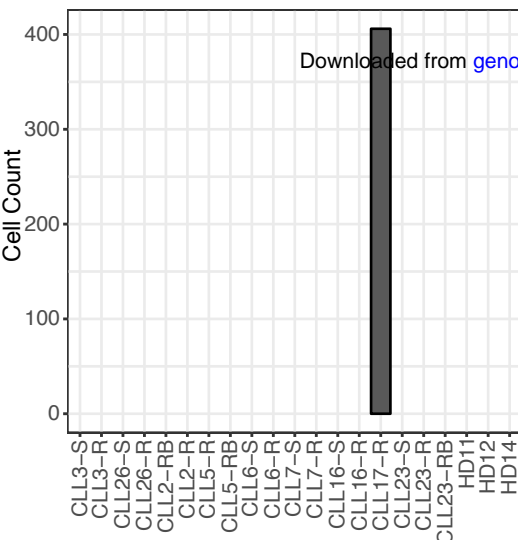




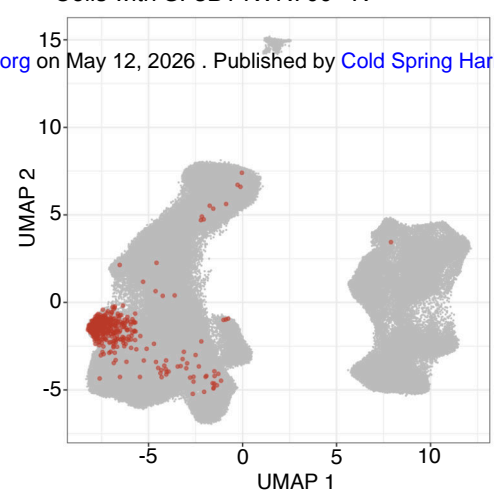




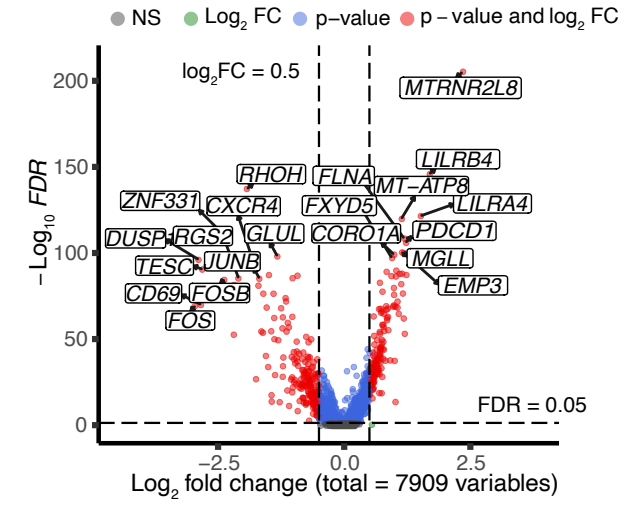
A Cells with > 2 UMIs of SF3B1 KVR700**R



B Cells with SF3B1 KVR700**R



C DEGs between SF3B1 KVR700**R and WT CLL cells



D DEGs between SF3B1 KVR700**R and WT CLL from CLL17

