

1 Editor summary: The authors introduce two cellular barcoding tools: CellBarcode, for extracting and filtering
 2 diverse DNA barcodes from bulk and single cell sequencing data; and CellBarcodeSim, for simulating
 3 barcoding experiments, thus enabling the investigation of the impact of biological/technical factors on
 4 barcode detection.

5

6 Peer review information: Primary Handling Editor: Fernando Chirigati, in collaboration with the Nature
 7 Computational Science team.

8

9 Peer review information: Nature Computational Science thanks Jennifer Adair, Mark Enstrom and the other,
 10 anonymous, reviewer(s) for their contribution to the peer review of this work.

11

12 **1. Supplementary Information:**

13 **A. PDF Files**

Item	Present?	Filename	A brief, numerical description of file contents.
		Whole original file name including extension. i.e.: Smith_SI.pdf. The extension must be .pdf	i.e.: <i>Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>
Supplementary Information	Yes	Supplementary_information_SunEtAl2024.pdf	Supplementary Tables 1-3, Supplementary Figures 1-25, Supplementary Vignettes 1 and 2, Supplementary Algorithm 1
Reporting Summary	Yes	NATCOMPUTSCI-23-0758Creporting.pdf	
Peer Review Information	No	<i>OFFICE USE ONLY</i>	

14

15 **B. Additional Supplementary Files**

Type	Number	Filename	Legend or Descriptive Caption
	Each type of file (Table, Video, etc.) should be numbered from 1 onwards. Multiple files of the same type should be listed in sequence, i.e.:	Whole original file name including extension. i.e.: <i>Smith_Supplementary_Video_1.mov</i>	Describe the contents of the file

	Supplementary Video 1, Supplementary Video 2, etc.		
Supplementary Data	Supplementary data 1	Supplementary_data_1_Igor_final_parms.txt	Parameters of the Bayesian network model inferred by IGoR from VDJ barcode data in mouse mammary gland tissue, then used to generate VDJ barcodes for the simulation study
Supplementary Data	Supplementary data 2	Supplementary_data_2_Igor_final_marginals.txt	Marginals of the Bayesian network model inferred by IGoR from VDJ barcode data in mouse mammary gland tissue, then used to generate VDJ barcodes for the simulation study

16

17 3. Source Data

Parent Figure or Table	Filename	Data description
	Whole original file name including extension. i.e.: <i>Smith_SourceData_Fig1.xls</i> , or <i>Smith_Unmodified_Gels_Fig1.pdf</i>	i.e.: Unprocessed western Blots and/or gels, Statistical Source Data, etc.
Source Data Fig. 1	n/a	
Source Data Fig. 2	Figure_2.xlsx	Statistical Source Data
Source Data Fig. 3	Figure_3.xlsx	Statistical Source Data
Source Data Fig. 4	Figure_4.xlsx	Statistical Source Data
Source Data Fig. 5	Figure_5.xlsx	Statistical Source Data
Source Data Fig. 6	Figure_6.xlsx	Statistical Source Data

18

19

20 Extracting, filtering and simulating 21 cellular barcodes using CellBarcode 22 tools

23 Author list

24 Wenjie Sun¹, Meghan Perkins², Mathilde Huyghe², Marisa M. Faraldo², Silvia Fre², Leïla Perié^{1,#},
25 Anne-Marie Lyne^{1,3,4,#}

26

27 Affiliations

28 1.Institut Curie, Université PSL, Sorbonne Université, CNRS UMR168, Laboratoire Physico
29 Chimie Curie, 75005, Paris, France

30 2.Institut Curie, Laboratory of Genetics and Developmental Biology, PSL Research University,
31 INSERM U934, CNRS UMR3215, 75005, Paris, France

32 3.INSERM U900, Paris, France

33 4.MINES ParisTech, CBIO-Centre for Computational Biology, PSL Research University, 75006
34 Paris, France

35 # These authors jointly supervised this work

36 Corresponding authors: sunwjie@gmail.com, anne-marie.lyne@curie.fr and leila.perie@curie.fr

37

38 Abstract

39 Identifying true DNA cellular barcodes amongst polymerase chain reaction (PCR) and sequencing
40 errors is challenging. Current tools are restricted in the diversity of barcode types supported or
41 the analysis strategies implemented. As such, there is a need for more versatile and efficient tools
42 for barcode extraction, as well as for tools to investigate which factors impact barcode detection
43 and which filtering strategies to best apply. Here, we introduce the package CellBarcode and its
44 barcode simulation kit, CellBarcodeSim, that allows efficient and versatile barcode extraction and
45 filtering for a range of barcode types from bulk or single-cell sequencing data using a variety of
46 filtering strategies. Using the barcode simulation kit and biological data, we explored the technical

47 and biological factors influencing barcode identification and provided a decision tree on how to
48 optimize barcode identification for different barcode settings. We believe that CellBarcode and
49 CellBarcodeSim have the capability to enhance the reproducibility and interpretation of barcode
50 results across studies.

51

52 Main text

53 Introduction

54 DNA cellular barcoding is a high-throughput approach widely used to follow lineage^{1,2} in different
55 fields such as hematopoiesis, development³⁻⁵, cancer⁶⁻⁹ and infection dynamics¹⁰. It uses unique
56 and heritable DNA sequence incorporated into the genome of an ancestor cell then detected via
57 sequencing in its progenies.

58

59 In the earliest approaches, progenitor cells were prospectively transduced ex vivo with libraries
60 of fixed-length oligonucleotides¹¹. More recently, to avoid extraction and reimplantation of
61 progenitor cells, in vivo recombining genetic cassettes have been incorporated in transgenic
62 organisms. Many innovative approaches have produce these in situ genetic labels¹²⁻¹⁶, with the
63 majority detected via short-read sequencing. Barcodes are now detected with single cell RNA-
64 sequencing (scRNA-seq)¹⁴⁻¹⁷, coupling lineage with fine-grained phenotyping.

65

66 DNA barcodes detected via Next Generation Sequencing (NGS) are subject to various sources
67 of error, resulting in the identification of spurious barcodes. All barcode types are affected by PCR
68 error/bias¹⁸ and sequencing error; in situ barcodes suffer additionally from the inability to control
69 the distance between barcodes^{19,20}. Biological factors such as the number of barcodes and clone
70 size can impact barcode detection but have rarely been investigated²¹. To extract and identify true
71 from spurious barcodes, many different bioinformatic filtering strategies have been proposed.
72 However, little comparison of the various strategies has been published and most publications
73 use their own “in house” processing pipelines. This is problematic in terms of interpretation of
74 results across studies and reproducibility. Both guidelines on how filtering strategies and their
75 parameterization impact barcode quantification as well as broadly applicable tools are required²².

76

77 Beside tools for visualization and data exploration²³⁻²⁵ three tools have been developed to extract
78 DNA barcodes from NGS data: genBaRcode²⁶, Bartender²⁷ and CellTagR²⁸. Whilst each has

79 demonstrated utility, they are either restricted in the diversity of barcode types supported
80 (CellTagR, genBaRcode) or the analysis strategies implemented (all of the above). No tools
81 provide a framework to simulate barcode experiments and investigate the technical and biological
82 factors impacting barcode detection. There is a need for more versatile tools to extract, identify
83 and simulate barcodes.

84

85 To address these issues, we developed two tools: CellBarcode, an R Bioconductor package for
86 barcode extraction and filtering, and CellBarcodeSim, a barcode simulation kit which faithfully
87 reproduces barcoding experiments. We demonstrate, using simulated and experimental datasets,
88 that CellBarcode allows users to implement various filtering strategies for bulk or single cell
89 datasets. Using CellBarcodeSim to simulate barcoding experiments, we investigated potential
90 technical and biological factors impacting the reliability of barcode identification, confirmed with
91 experimental datasets. We recapitulated our results into a decision tree to guide researchers on
92 which filtering strategy is most appropriate for their setting. Overall, we present efficient and
93 versatile tools to extract and identify barcodes from errors, and provide advice on how best to
94 analyze barcoding experiments in a range of biological situations.

95

96 Results

97

98 CellBarcode package

99

100 We developed the CellBarcode R package, which provides a toolkit for barcode pre-processing,
101 including steps from generating the FASTQ quality control (QC) information to exporting the data
102 into a read count matrix (Figure 1A). Using the read QC & filtering functions of CellBarcode, users
103 can check sequencing quality, remove low-quality sequences, and get an overview of read
104 diversity. Barcodes can then be extracted from the FASTQ or BAM file by defining a regular
105 expression matching the structure of the lineage barcode and its surrounding flanking sequence
106 (see Supplementary vignette 1 for examples and a detailed description of this process); both
107 fixed-length and variable-length barcodes can be extracted, and mismatches in the flanking
108 regions are allowed (bulk analysis only). Once the raw barcodes have been extracted, filtering
109 functions can remove spurious barcode sequences using commonly applied strategies. In
110 addition, the package provides functions for visualizing the barcode read count distribution per
111 sample and across replicates (Figure 1B and 1C).

112

113 The four main filtering strategies generally applied to barcoded data are implemented in
114 CellBarcode (Figure 1D): 1) Reference filtering: barcodes not matching with the reference list are
115 eliminated. The reference list is either generated by sequencing the viral barcode libraries⁵ or
116 enumerating all possible barcodes using knowledge of barcode structure¹⁹; 2) Threshold filtering:
117 barcodes are retained if their read number (depth) surpasses a specified threshold⁵. CellBarcode
118 has a manual or an automatic threshold option (Methods section ‘Barcode Filtering’); 3). Cluster
119 filtering: barcodes that have an edit distance smaller than a specified threshold to a more
120 abundant barcode are eliminated²⁹; 4) unique molecular identifier (UMI) filtering: if UMIs are
121 added to DNA molecules during library preparation, several optional filtering steps can be applied,
122 including extracting the most abundant barcode per UMI and threshold filters on the read count
123 per UMI or UMI count per barcode. These 4 filtering strategies can be used individually or in
124 combination, and we later advice on when to apply each strategy using simulated data with
125 CellBarcodeSim. See Supplementary Vignettes 1 and 2 for examples of all major use cases.

126

127 In summary, CellBarcode is a versatile and open-source tool that works on all major operating
128 systems and is capable of analyzing a wide variety of DNA barcode types with commonly applied
129 filtering strategies. The key assets of CellBarcode are its speed, the ability to deal with UMI data
130 and the extraction of barcodes from scRNA-seq data (Supplementary Table 1). Efficient C++ code
131 accelerates heavy tasks compared to other packages; barcode extraction and cluster filtering are
132 20 and 70 times faster than using genBaRcode (Supplementary Figure 1).

133

134 Comparing barcode filtering strategies using CellBarcodeSim

135

136 The CellBarcode package provides a variety of functions for barcode filtering, but choosing a
137 filtering strategy and its parameterization in a given experimental setting is challenging. With this
138 in mind, we developed a barcode simulation toolkit, called CellBarcodeSim which produces in
139 silico barcoding data mimicking bulk DNA-seq experimental situations by varying a number of
140 technical and biological factors. CellBarcodeSim covers production of a barcode library, cell
141 barcode labeling and clonal expansion, construction of full sequencing reads including flanking
142 sequences and UMIs when desired, and finally PCR amplification and sequencing with the
143 inclusion of error (Figure 2A, Methods section ‘DNA cellular barcode sequencing simulations’). In
144 total, CellBarcodeSim provides 10 configurable parameters for non-UMI and 13 for UMI
145 sequencing libraries (Figure 2A). Tens of thousands of clones can be simulated on a standard

146 laptop (16 Gb random-access memory), covering most experimental situations. Two types of
147 barcode libraries can be simulated with CellBarcodeSim (Methods section ‘DNA cellular barcode
148 sequencing simulations’) while other types of barcodes can be uploaded as a list. Comparing the
149 known barcodes from simulation with the output of CellBarcode can guide users in their choice of
150 filtering strategy and its parameterization. Overall, CellBarcodeSim simulates barcoding
151 experiments varying multiple technical and biological factors.

152

153 Before exploring how different parameters impact barcode identification across filtering strategies,
154 we first checked that CellBarcodeSim could reproduce the expected output of a barcoding
155 experiment. We simulated two experimental datasets: lentiviral fixed-length 20bp barcodes
156 recovered from myeloid cells³⁰; and a variable, diversity, joining (VDJ)-barcoded dataset with
157 UMIs recovered from Mouse Embryonic Fibroblast (MEF) cells²⁰ (Methods section ‘Acquisition,
158 analysis and simulation of experimental data’). We showed that CellBarcodeSim outputs the same
159 read structure and similar proportion of reads matching the regular expression as the
160 experimental data (Figure 2B and 2C), with high Pearson correlation between the proportion of
161 the most abundant base at each sequencing cycle between the simulated and experimental data
162 (Figure 2B and 2C).

163

164 Next, to investigate the key factors impacting barcode identification for different filtering strategies,
165 we first designed a default scenario for non-UMI barcode libraries (Methods section ‘DNA cellular
166 barcode sequencing simulations’ and Supplementary Table 2) and then 25 alternative scenarios
167 varying key biological and experimental parameters (Supplementary Table 2). After randomly
168 simulating each scenario 30 times, we applied 4 different filtering strategies (read count
169 thresholding, reference library, clustering and UMI filtering). To evaluate the filtering performance,
170 for each simulation we computed barcode recall (the proportion of true barcodes found in the
171 output) and precision (the proportion of output barcodes which are true) using the known ground
172 truth. We then computed the area under the precision-recall curve (P-R AUC) across a range of
173 thresholds (Supplementary Figure 2) to indicate how well filtering methods separate true from
174 spurious barcodes regardless of threshold.

175

176 We first consider read count threshold filtering. In all scenarios, there is an overlap between the
177 read count distributions of error and true barcodes combined across simulations (Supplementary
178 Figure 3), therefore it is impossible to choose a read threshold to perfectly separate true from
179 spurious barcodes. Using a read threshold involves a trade-off between the recall and precision

180 of barcode detection, with higher threshold removing more spurious barcodes but also more true
181 barcodes (Figure 3A and 3B). Surprisingly, the factor which had the largest impact on P-R AUC
182 was one of the biological factors: the standard deviation (SD) of the log clone size (where log
183 denotes the natural logarithm), with smaller clone size variation displaying larger P-R AUC (Figure
184 3C, Supplementary Figure 4 and Supplementary Figure 5A). When log clone size SD was 1, the
185 P-R AUC reached 1 regardless of other factors including barcode type or mean clone size (Figure
186 3C, Supplementary Figure 4 and Supplementary Figure 5A). Comparing precision and recall for
187 different thresholds, we observed the expected trend of increased recall but decreased precision
188 as the threshold became less stringent (Figure 3A, Figure 3B, Supplementary Figure 6 and
189 Supplementary Figure 5B). When there is high variability in the number of cells labeled by each
190 barcode (log clone size $SD \geq 2$), recall needs to be compromised to avoid calling spurious
191 barcodes. This leads to a significant loss of true barcodes, predominantly affecting barcodes of
192 small clones that have similar read count to error barcodes derived from much larger clones
193 (Figure 3A, Figure 3B, Supplementary Figure 6 and Supplementary Figure 5B). This loss of true
194 barcodes can preclude robust statistical analysis downstream (Supplementary Figure 5C).

195

196 To validate the finding about the impact of clone size SD, we used an unpublished dataset in
197 which Cas9-expressing mice intestinal organoids were infected with libraries of gRNAs designed
198 to knock out specific genes (Methods section 'CRISPR gRNA dataset'). Whilst not a standard
199 barcode, each specific knock-out acts as a clonal label and can be extracted by CellBarcode using
200 a regular expression targeting the constant primer region. Two time points were analyzed, 24
201 hours and 7 days, with clone size variation increasing over time due to fitness effects of the
202 gRNAs. Using CellBarcodeSim to simulate the experiment, we successfully reproduced the
203 percentage of barcode-containing reads, and observed a change in the read count distribution,
204 from bimodal with true and spurious barcode counts mostly separated at low clone size SD, to
205 unimodal with more overlap in true and spurious barcode counts at higher clone size SD
206 (Supplementary Figure 7, top row). These same trends were observed in the experimental data
207 (Supplementary Figure 7, bottom row). To verify the finding that number of PCR cycles has limited
208 impact on barcode recall (Supplementary Figure 8), we used published data of mixes of 7 MEF
209 cell lines that each contain a unique known VDJ barcode²⁰. Across the mixes, the total number of
210 initiating cells was reduced and the number of PCR cycles correspondingly increased to produce
211 a constant PCR product concentration, with the clone size ratios kept constant. Irrespective of the
212 number of PCR cycles, CellBarcode identified the 7 known barcodes in each mix (with one
213 spurious barcode at +4 PCR cycles) (Supplementary Figure 8). Using CellBarcodeSim with

214 matched parameters and varying the number of PCR cycles, we reproduced the separation of true
215 and spurious barcode counts and the lack of change in the sequence frequency distribution
216 (Supplementary Figure 8). Using two experimental datasets, we therefore demonstrated that
217 CellBarcodeSim can simulate real scenarios. Our simulation results of the large impact of the clone
218 size-SD and the limited impact of PCR cycle number on barcode identification were supported by
219 these experimental data. Regarding filtering, we showed that the read count thresholding strategy
220 is suboptimal at best, except for systems in which the clones have a similar number of cells. Some
221 biological systems have been shown to differ in their proliferation capacities³¹, but for most of
222 them this information is unknown. CellBarcodeSim is therefore a useful tool to simulate different
223 scenarios, guiding researchers on the impact of thresholds on barcode identification and aiding
224 in the interpretation of results.

225

226 An alternative strategy for barcode filtering is to match the extracted barcodes to a reference
227 library when available. Using this approach for fixed-length barcodes, the distributions of true
228 barcode read counts overlap less with those of spurious barcodes (Supplementary Figure 9), and
229 true barcode P-R AUC was substantially improved, with most scenarios having a P-R AUC of 1
230 (Supplementary Figure 10), as suggested before⁵. We applied read count thresholding here after
231 reference filtering to compute the P-R AUC enabling scenario comparison, although its use is
232 optional. We note that read count threshold filtering is used to call true barcodes in the generation
233 of the reference library itself, and even though these plasmid libraries have more homogenous
234 barcode abundances than most biological experiments, the reference library suffers from the
235 threshold-related problems described above and by others²¹. For variable-length barcodes such
236 as VDJ barcodes, a reference library can be generated by simulating all possible combinations.
237 Using this list had limited improvement in P-R AUC (Figure 3D, Supplementary Figure 10) due to
238 the small edit distance between some barcodes (many with edit distance < 3, Supplementary
239 Figure 11A). Spurious sequences created by PCR or sequencing error can have the same
240 sequence as a barcode in the reference library (Figure 3F) and are not filtered out, impacting the
241 precision (Supplementary Figure 12). Overall, these results show that a reference library is a
242 useful approach for fixed-length barcodes designed to have edit distances larger than 3, but is
243 not useful for variable-length barcodes such as VDJ barcodes where the edit distance cannot be
244 controlled.

245

246 Several studies have advocated cluster filtering to identify true barcodes^{21,26,27}. With clustering,
247 true barcodes are identified by comparing barcode sequences, usually with the assumptions that

248 barcodes separated by very short edit distances are the result of PCR/sequencing errors and that
249 the most abundant barcode in the cluster is the true barcode^{21,26,27}. We used CellBarcodeSim to
250 evaluate how cluster filtering performs compared to other filtering strategies. Cluster filtering
251 improved the P-R AUC of random barcodes compared to threshold filtering alone (Figure 3E) and
252 performed as well as reference library filtering (Figures 3D and E), implying that it is the method
253 of choice for the generation of a reference library, as previously suggested^{18,21}. For variable length
254 barcodes like VDJ barcodes, clustering performed worse or similar to threshold or reference
255 library filtering (Figure 3E and Supplementary Figure 13) due to low recall (Supplementary Figure
256 14), although the true barcode read counts overlap less with those of the spurious ones
257 (Supplementary Figure 14). This is linked to the short edit distance of some in-situ barcodes,
258 which are not PCR/sequencing errors as assumed by cluster filtering (Figure 3F and
259 Supplementary Figure 11A). We previously developed a sequencing library preparation protocol
260 for VDJ barcodes with UMIs²⁰. We hypothesized that the addition of UMIs will improve the
261 identification of true barcodes using cluster filtering. To test this hypothesis, we simulated VDJ
262 barcode sequencing with UMIs for high clone size variation samples, which we identified as the
263 most difficult scenario in which to apply this filtering (Supplementary Table 3). We observed that
264 incorporating UMI information significantly improved the P-R AUC for samples with large clone
265 size variation (Figure 3G), supporting the hypothesis that the addition of UMIs helps true barcode
266 identification by cluster filtering for barcodes with low edit distance, such as VDJ barcodes.
267 Overall, these results show that cluster filtering is an efficient method to identify barcodes in
268 systems with large edit distance such as viral barcodes^{18,32}. It is the method of choice if one had
269 no reference library or to make a reference library for such barcodes^{18,21}.

270

271 We summarized the findings of our comprehensive comparison in a decision tree to guide
272 researchers on which strategy to apply to their data (Figure 3H). In summary, our advice is: use
273 reference library or cluster filtering if the barcoding system has a large edit distance
274 (approximately ≥ 3); otherwise if the barcode clone size variation is small, a read threshold would
275 work. If the barcode clone size variation is large and the barcode system has small edit distance,
276 either UMIs need to be used or a stringent read count threshold implemented sacrificing true
277 barcodes with low read count.

278

279 Reference and cluster filtering of lentiviral barcodes

280

281 To compare cluster and reference library filtering on biological data, we used CellBarcode to
282 analyze paired technical replicates of 13,564 myeloid cells labelled with a random fixed-length
283 barcode library³⁰. Consistent with simulated random barcodes (Supplementary Figure 11B) it
284 displayed a high edit distance (Supplementary Figure 11C). First, we used CellBarcode to check
285 the quality of the FASTQ file, plotting the base percentage and quality in each sequencing cycle
286 (Figures 4A and B). We successfully extracted and quantified the barcodes using CellBarcode as
287 shown by the correlation with those in the original paper (Figure 4C). Our results are also
288 consistent with genBaRcode (Supplementary Figure 15A) and Bartender analysis
289 (Supplementary Figure 15B), although we observe considerably more noise in the Bartender data,
290 because it has fewer filtering steps implemented.

291
292 According to our decision tree, the methods to use for high edit distance barcodes are reference
293 library or cluster filtering. We therefore extracted barcodes using either no filtering, reference
294 library or cluster filtering and compared barcode cell count detected in technical repeats after
295 normalizing read counts by total cell number (Figure 4D). In biological data as the identity of the
296 true barcodes is unknown, we used the reference library provided in³⁰. Without filtering, many
297 barcodes not present in the reference library overlapped in read count distribution with those in
298 the reference library, agreeing with our simulation results that read threshold filtering decrease
299 the recall to ensure precision (Figure 4D). Cluster filtering removed most of the barcodes absent
300 from the reference library, leaving only one spurious sequence present in one cell, while keeping
301 all the true barcodes with more than one cell (Figure 4E). This confirms our simulation finding that
302 cluster filtering can have the same efficacy as reference library filtering using barcodes with high
303 edit distances.

304 305 Read threshold filtering of in situ barcodes

306
307 Variable-length barcodes like VDJ barcodes are the most challenging to identify in noisy data due
308 to the short edit distance barcodes generated. To explore if our CellBarcode simulation results
309 would hold in experimental variable-length barcode data, we made use of our unpublished in vivo
310 VDJ barcode data from mouse mammary glands, for which we have both UMI and non-UMI data
311 from the same sample (Figure 5A and B). Using the known read structures of the two sequencing
312 libraries (Figure 5B), we extracted the barcodes and applied automatic read threshold filtering
313 and UMI filtering to the non-UMI and UMI samples respectively (Figure 5C and Figure 5D),
314 illustrating the versatility of CellBarcode to extract barcodes from a variety of structures (Methods

315 section 'VDJ barcode mammary gland dataset). For different UMI read count thresholds, we
316 observed that the number of barcodes reached a plateau (Supplementary Figure 16A). At this
317 plateau, in one duplicate sample, we identified 80 barcodes in the non-UMI library, and 82
318 barcodes in the UMI library with 76 barcodes overlapping (87%) (Figure 5E).

319
320 In this data, the biggest clones had about 100 times higher read/UMI count compared to the
321 smallest clones, corresponding to a log clone size SD of 1, the lowest considered in our
322 simulations (Figures 5C and D). The clone sizes in the UMI and non-UMI libraries after threshold
323 filtering (normalized reads or UMI count) correlated very well (Figure 5F), with most of the
324 inconsistent barcodes being small clones. This result supports our simulation conclusion that
325 automatic read thresholding performs well in experimental settings with small clone size variation.
326 We observed more spurious barcodes in both UMI and non-UMI results from Bartender (see
327 Supplementary Figure 16B and C), indicating the importance of read or UMI read count thresholds
328 which are not implemented in Bartender.

329

330 Using CellBarcode to analyze scRNA-seq data

331

332 Finally, we designed CellBarcode to extract and identify lineage barcodes from single cell omics
333 data. To this end, CellBarcode is equipped with functions to process barcodes from the most
334 popular technologies such as 10x Genomics or Smart-seq (Figure 6A). In this section, we use the
335 term 'cell barcode' to refer to the unique barcode labeling each cell from the single cell sequencing
336 protocol, and 'lineage barcode' to refer to the barcode added during a lineage tracing experiment.
337 Input to CellBarcode is flexible, allowing FASTQ and BAM/SAM files, either one file for all cells
338 (as for 10x Genomics scRNA-seq) or one file per cell (as for Smart-seq2), and BAM/SAM files
339 pre-tagged with cell barcodes and UMIs, such as those output by the 10x Genomics software
340 CellRanger. We illustrate the use of CellBarcode on scRNA-seq data but it applies to many types
341 of lineage barcoded single cell omics data, such as scATAC-seq^{33,34}. The potential (but optional)
342 filters include 1). Extract dominant barcode per UMI, 2). Filter UMIs using a read count threshold
343 and 3). Filter lineage barcodes using a UMI count threshold (Figure 6B). The user must choose
344 various thresholds, and here we distinguish two experimental scenarios from published data: 1.
345 a unique lineage barcode per cell, such as low concentration lentivirus infection¹⁷ or heterozygous
346 inducible VDJ barcode³⁵, and 2. multi-barcodes per cell, for example, high concentration lentiviral
347 infection such as the CellTag barcode system³⁶.

348

349 To compare the performance of CellBarcode to that of CellTagR, a dedicated package for analysis
350 of barcoded scRNA-seq data, we replicated the CellTagR demo analysis pipeline
351 (<https://github.com/morris-lab/CellTagR>) with CellBarcode on the multi-barcode per cell data
352 from³⁶. Applying the same steps and parameters (Methods section 'CellTag barcode scRNA-seq
353 dataset'), CellBarcode obtained similar results to CellTagR (Supplementary Figure 17A and B)
354 with 20% less runtime (Supplementary Figure 17C and D). CellTagR only supports the extraction
355 of CellTag barcodes, whereas, to illustrate the versatility of CellBarcode, we extracted variable-
356 length VDJ barcodes from scRNA-seq data from Cosgrove et al (2023)³⁵ and obtained similar
357 barcodes and quantification to the original paper (Supplementary Figure 18).

358

359 To illustrate how CellBarcode can help users select the different filtering thresholds, we counted
360 the number of lineage barcodes retrieved per cell for various types of filtering in VDJ barcoding
361 data from Cosgrove et al (2023)³⁵. Due to the introduction of one VDJ cassette in one allele of the
362 mouse genome, each cell in this dataset has only one lineage barcode. We observed a trade-off
363 between the accuracy of lineage barcode retrieval (i.e. the proportion of cells with one unique
364 lineage barcode) and the total number of lineage-barcoded cells retained for analysis. We first
365 filtered to take the dominant lineage barcode per UMI, as the combination of high-diversity cell
366 and UMI barcodes for each read can be assumed unique, which dramatically reduced the number
367 of barcodes per cell in comparison to the raw data (Figure 6C). Using different minimum read-
368 count-per-UMI thresholds, we found that the number of barcodes per cell was easily restricted to
369 a maximum of 2 with a threshold of 2 (Figure 6D). Increasing the read-count-per-UMI threshold
370 further resulted in the loss of many cells for analysis (Figure 6D). Complementing the read-count-
371 per-UMI filtering with a UMI-count-per-barcode filter of 2, we obtained one identifiable lineage
372 barcode per cell (Figure 6E). These thresholds will depend on each specific dataset, for example,
373 with low sequencing depth, even without read-count-per-UMI or UMI-count-per-cell filtering, most
374 cells have one unique lineage barcode as observed in the Marsolier et al. (2022)¹⁷ dataset
375 (Supplementary Figure 19).

376

377 To conclude, in addition to an improvement in run time, CellBarcode can extract and identify
378 lineage barcodes in scRNA-seq data from many different barcode designs due to its flexible use
379 of regular expressions. Moreover, CellBarcode implements several filtering strategies to identify
380 true from spurious lineage barcodes in single cell data, and produces figures helping the user
381 choose a strategy and its parameterization.

382 Discussion

383 In this paper, we presented CellBarcode, a versatile R package for analysis of barcoding data,
384 and CellBarcodeSim, a pipeline to simulate barcoding experiments. Whilst we designed the
385 simulation tool to test and parameterize filtering approaches for barcode identification, it can be
386 employed in a similar vein for experimental design; for example, users can investigate the impact
387 of different barcode lengths, UMI or no-UMI libraries and sequencing depths in their biological
388 scenario. We would highlight, however, that this is complicated by the combination of unknown
389 biological factors and final filtering approach.

390
391 Beltman et al. (2016)²¹ suggested not to use cluster filtering as it can result in the removal of true
392 barcodes. However, both our simulations and tests on real data show that cluster filtering
393 performs well when the barcode edit distance is large enough (≥ 3 in our simulations) compared
394 to realistic low levels of PCR/sequencing error. We would therefore refine the statement from
395 Beltman et al (2016)²¹ to add that cluster filtering can be successfully used when the edit distance
396 is sufficiently high, even in the case of high clone size variation.

397
398 We modeled clone size using a log-normal distribution based on our analysis of t-cell receptor
399 clones (Supplementary Figure 20), and whilst users of CellBarcodeSim can also opt for a power
400 law distribution, we hope to add more detailed models in future versions of the tool (such as one
401 based on Radtke et al. 2023³⁷). Indeed, in most systems the clone size distribution is unknown;
402 in this case CellBarcodeSim can be used to investigate the impact of filtering strategies on
403 barcode identification under different assumptions and can aid users in their biological
404 interpretation. Further simulation work is also required to test the impact of filtering on barcode
405 quantification.

406
407 CellBarcodeSim makes many other assumptions about the processes involved to simulate
408 barcoding data. Barcode library production is modeled with simple distributions rather than
409 separately modeling the stages of transfection, growth and sampling. The fixed-length Hamming³⁸
410 barcodes simulated using the DNABarcodes package are filtered to remove many sources of
411 error problematic for PCR, such as barcodes containing triplets or with GC bias. The PCR
412 simulation assumes the amount of starting material is large enough to ignore contamination and
413 doesn't model factors such as non-specific hybridizations. Indeed, we do not expect our simulation
414 to quantitatively model all possible effects of the complex PCR process. Researchers interested

415 in specific sources of error, such as those introduced during barcode library preparation, or using
416 a non-standard protocol where the PCR primer does not target the constant flanking region, would
417 need to adapt the simulation.

418

419 CellBarcodeSim calls external tools such as ART NGS read simulator³⁹, DNABarcodes R
420 package to simulate fixed-length barcodes⁴⁰ and IGoR to simulate VDJ barcodes⁴¹, which could
421 be a concern in terms of longevity. ART is a mature and heavily used tool with no updates required
422 and containing pre-built error profiles for all the major sequencers. The packages simulating
423 barcodes are less mature and barcode type specific, but CellBarcodeSim can be easily updated
424 allowing other tools to feed in.

425

426 Methods

427

428 Ethics statement

429 All studies and procedures involving animals were in accordance with the recommendations of
430 the European Community (2010/63/UE) for the Protection of Vertebrate Animals used for
431 Experimental and other Scientific Purposes. Approval was provided by the ethics committee of
432 the French Ministry of Research (reference APAFIS #34364-202112151422480). We comply with
433 internationally established principles of replacement, reduction, and refinement in accordance
434 with the Guide for the Care and Use of Laboratory Animals (NRC 2011). Husbandry, supply of
435 animals, as well as maintenance and care in the Animal Facility of Institut Curie (facility license
436 #C75-05-18) before and during experiments fully satisfied the animal's needs and welfare.
437 Mouse breeding was in a specific pathogen-free animal facility and animals were co-housed with
438 housing conditions using a 12 light/12 dark cycle, temperature between 20 and 24 °C and average
439 humidity rate between 40% and 70%.

440 DNA cellular barcode sequencing simulations

441 We simulated the DNA cellular barcode sequencing data using CellBarcodeSim (version 1.0) with
442 5 steps: 1). Lineage barcode simulation, 2). barcode labeling 3). clonal expansion, 4). PCR
443 amplification, and 5). sequencing.

444

445 Lineage barcode simulation

446 Two types of barcode libraries can be simulated with CellBarcodeSim ("random barcodes" with
447 uniform probability and fixed-length, and "Hamming barcodes" with uniform probability, fixed-
448 length and a minimum Hamming distance between sequences) while other types of barcodes can
449 be uploaded as a list. In addition, three libraries were simulated and uploaded in the package: 14
450 base pair (bp) random barcodes, 14bp Hamming barcodes with minimum distance 3 simulated
451 using DNABarcode⁴⁰, and variable length VDJ barcodes²⁰ simulated using an external package
452 IGoR.

453 For the simulation study, a list of possible barcodes was simulated for three types of barcode and
454 barcodes were randomly sampled from this list to label cells. The fixed-length uniform-probability
455 'random barcodes' were generated with `stri_rand_strings` from `stringi` package. To generate
456 'Hamming barcodes' with a minimum Hamming distance of 3, we used the `create.dnabarcodes`
457 function from the `DNABarcodes` package⁴⁰. The barcode length can be defined by the user. In
458 this simulation study, we tested 14 or 10 base-pair. Lastly, for the variable length 'VDJ barcodes'²⁰,
459 a list of 1×10^7 VDJ barcodes to sample from was generated using IGoR⁴¹. To ensure the simulated
460 VDJ barcodes resemble those produced in vivo, the parameters of the Bayesian network model
461 used to generate the barcode space were inferred using IGoR from the VDJ barcode sequencing
462 data in mammary gland tissue (Supplementary Data 1 & 2). Among the simulated sequences,
463 there are 1.4×10^5 unique barcode sequences with different frequencies. To simulate the noise
464 during library preparation for random or Hamming barcodes, CellBarcodeSim can simulate
465 normal, log-normal or exponential distributions, or the user can simulate according to their own
466 uploaded empirical distribution.

467

468 Barcode labeling simulation

469 We randomly sampled the barcode lists simulated in the previous step for the corresponding
470 barcode type. We simulated different samples with different total barcode numbers. Each barcode
471 labels one initial cell in the simulation, and those barcode sequences were used as the true
472 barcodes in later precision & recall analysis. We tested scenarios with 300-30,000 initiating cells,
473 but as we found the sequence count distributions to be very similar, as well as the impact of
474 various factors on the precision and recall, we chose values of 150, 300, 600 and 1,200 for the
475 repeat simulations, corresponding to the number of barcodes in most published work.

476

477 Clonal expansion simulation

478 We used a log-normal distribution to simulate the final clone sizes of the initially labeled cells. The
479 parameters of the reference distribution are log mean 1.2 and log standard deviation (SD) 2, which
480 were chosen based on the experimentally-derived murine naive CD8 TCR beta chain sequence
481 clone size distribution described in Desponds, Mora and Walczak 2016 (Supplementary Figure
482 20)⁴². Observing log clone size SDs of ~1 in our VDJ barcoded mammary gland data, ~2.5 in
483 Eisele et al (2022)³⁰ and ~2.5-3 in Adair et al (2020)⁴³, we define alternative scenarios of log clone
484 size SD 1 and 3. We used the rlnorm function in R 4.2.1⁴⁴ to generate random numbers and the
485 clone size of each barcode clone was defined by rounding up the nearest integer of the
486 corresponding random number. The CellBarcodeSim tool also offers the power-law clone size
487 distribution.

488 We note that when the clone size follows a log-normal distribution, the ratio of the 99th quantile,
489 $Q(0.99)$, divided by the 1st quantile, $Q(0.01)$, depends only on the log standard deviation and not
490 on the log mean (Supplementary Figure 21), which is explained by the following equations:

491
$$Q(q) = e^{\mu + \sigma \cdot \Phi^{-1}(q)} \quad (1)$$

492 where $\Phi^{-1}(q)$ is the q -th quantile of the standard normal distribution with mean,
493 μ , and standard deviation, σ .

494 The ratio of the 99th quantile to the 1st quantile:

495
$$Q(0.99)/Q(0.01) = e^{\mu + \sigma \cdot \Phi^{-1}(0.99)} / e^{\mu + \sigma \cdot \Phi^{-1}(0.01)} = e^{\{\sigma \cdot (\Phi^{-1}(0.99) - \Phi^{-1}(0.01))\}} \quad (2)$$

496 Therefore, we can use the range of empirical clone sizes as a quick estimation of the log standard
497 deviation.

498

499 PCR expansion simulation

500 The PCR simulation was written in C++ and assumes exponential amplification with efficiency of
501 0.703⁴⁵ and error rate of 1×10^{-5} for Taq enzyme, 1×10^{-6} for Phusion enzyme, and 1×10^{-7} for Q5
502 enzyme. Since PCR mutations are rare events, it is unlikely to have more than one mutation per
503 sequence molecule per PCR cycle, and substitution errors are the dominant PCR error type⁴⁶.
504 We therefore only allow a maximum of one base substitution per PCR cycle. In the simulation, we
505 replicated the barcode DNA sequence in-silico with the probability of the amplification efficiency,

506 rounding to the nearest natural number, and randomly mutated the base of the newly synthesized
507 sequence with the PCR error rate. To reduce the memory usage, as most of the barcodes have
508 the same sequence due to the low PCR error, we stored barcode sequences in a frequency table
509 of barcode sequences and frequencies. For the new PCR products, the mutant molecular
510 abundance was estimated by multiplying each sequence frequency by the error ratio, considering
511 the sequence length. The value was rounded to the nearest integer. Then uniform random
512 numbers were generated to decide the mutation position and substitution base-pair. The
513 sequence frequency table was updated by integrating the mutant sequence. If using UMIs,
514 investigators can select the number of pre-UMI PCR cycles (in which the UMI sequence will not
515 accumulate PCR errors) and the number of post-UMI PCR cycles (when the UMI sequence will
516 accumulate PCR errors). Since the PCR primer region is unlikely to have a PCR mutation and
517 this generally corresponds to the barcode flanking regions, by default, the flanking sequence is
518 added after the PCR simulation, matching the sequence to the experimental case when
519 applicable. However, investigators have the option to include the flanking region in the PCR
520 simulation by appending the fixed flanking regions to the barcodes when simulating the barcode
521 library (see Supplementary Vignette 1 for more detail).

522

523 Sequencing simulation

524 Sequencing simulation was conducted using the ART (version 2016-06-05) command line tool (a
525 next-generation sequencing reads simulator), which supports base substitution, insertions, and
526 deletions³⁹. The ART-integrated MiSeq V1 and HiSeq2000 read error profiles (learnt empirically
527 from relevant training data³⁹) were used to generate single-end sequencing with 100 base pairs,
528 with other parameters as default. We describe the sequencing profiles used in Supplementary
529 Figure 22A and B, together with PCR error in Supplementary Figure 22C. When comparing the
530 barcode clone size distributions between different simulated datasets, we sample 10^5 sequencing
531 reads to make the distributions easier to compare.

532

533 Simulating VDJ barcoded data with high clone size variation and UMIs

534 We simulated VDJ barcode sequencing with UMIs for high clone size variation samples (details
535 of the parameters in Supplementary Table 3). With an expected sequencing depth of 50 reads
536 per UMI, we filtered out UMIs that have read < 10 (based on sensitivity analysis to identify when
537 the number of barcodes detected plateaus) and then varied the UMI count threshold to compute
538 the P-R AUC.

539 DNA cellular barcode pre-processing strategy evaluation

540

541 Evaluation of filtering strategies: precision, recall, AUC

542 In the simulation study, we evaluated filtering strategies using precision and recall. The precision
543 and recall are defined as:

$$544 \text{ Precision} = n_{\text{true}} / n_{\text{output}} \quad (3)$$

$$545 \text{ Recall} = n_{\text{true}} / n_{\text{input}} \quad (4)$$

546 where n_{input} is the number of barcodes used for labeling, n_{output} is the total number of barcodes in
547 the pre-processing output, n_{true} is the number of barcodes shared between the pre-processing
548 output and the barcodes used for labeling.

549

550 The precision and recall depend on the threshold used for barcode filtering. Precision-recall
551 curves were drawn using a range of read count thresholds (or UMI count in UMI cleaning case),
552 and the area under the curve (AUC) was calculated to evaluate the overall goodness of a filtering
553 strategy. The AUC is a way to evaluate the goodness of a method regardless of threshold and
554 was computed using the ROCR R package⁴⁷.

555

556 All boxplots depict 25, 50, and 75th percentiles in the box, 25th or 75th percentile minus or plus
557 1.5*IQR respectively for the whiskers and points show outliers beyond the whiskers.

558 Barcode filtering

559

560 We enabled four barcode filtering strategies in the CellBarcode package with `bc_cure_umi`,
561 `bc_cure_clustering`, `bc_cure_depth` and `bc_auto_cutoff` functions. They are 1). read count
562 thresholding filtering with `bc_cure_depth` function, 2). Reference library filtering, 3). cluster filtering
563 and 4). UMI filtering.

564

565 **Read count threshold filtering** excludes the barcodes with read counts under the threshold. The
566 automatic threshold function determines the threshold by applying 1-dimensional weighted k-
567 means clustering to the barcode read count distribution. It involves the following steps: 1).
568 Remove barcodes with count below the median (as there are generally many more spurious than
569 true barcodes). 2). Transform counts by $\log_2(x+1)$. 3). Apply 1-dimensional k-means clustering⁴⁸
570 to the transformed read counts with cluster number fixed at 2 and with weights of the transformed
571 count. 4). Use the boundary between the two clusters as the read count threshold.

572

573 In **reference library filtering**, only barcodes appearing in the barcode reference list are retained
574 in the final output, and all others are filtered out. In the simulations, the barcode reference library
575 was the barcode list generated in the “Lineage barcode simulation”.

576

577 For **cluster filtering**, we assumed that with a low error rate, spurious error barcodes should have
578 a much lower read number compared to their true “mother” sequences. We clustered barcodes
579 with similar sequences to identify potential “mother” and “daughter” sequence pairs. Then we
580 removed the “daughter” sequences, thus making it easier to identify true barcodes with small
581 clone size. We used the following clustering process for each sample: 1). Identify the most
582 abundant barcode based on read counts; 2). Compute the distance (Hamming distance or
583 Levenshtein distance) between the most abundant barcode and the other barcodes, starting from
584 the least abundant barcode; 3). If the distance between two barcodes is below a set threshold,
585 and the reads count fold change between them is above a set threshold, the less abundant
586 barcode is removed; 3). Iterate for each of the other barcodes in order of abundance. The process
587 is described by the pseudo code in Supplementary Algorithm 1.

588

589 **UMI filtering** takes advantage of the unique molecular identifier (UMI) sequence. The default in
590 CellBarcode is to assume UMIs are not unique in line with the findings of Venkataram et al⁴⁹
591 (although the reader has the option to assume the converse if they wish). We first counted the
592 number of reads for each UMI-barcode combination and then applied a read count threshold. The
593 remaining barcode abundances were quantified by summing the UMI count. We assume that the
594 probability of an error in both the UMI and its associated barcode sequence is very low, and so
595 we do not cluster similar UMIs. This may result in a slight overestimation of clone size if a UMI
596 sequence results from an error, but should not affect barcode identification.

597

598 Benchmarking CellBarcode and genBaRcode

599 In order to compare the output and run time of CellBarcode (version 1.7.1) and genBaRcode (with
600 version 1.2.6), we simulated a random barcode dataset using the method described above with
601 parameters 1). 300 cells induced, 2). Log-normal clone size distribution with log clone size SD of
602 2 and log clone size mean 1.2, 3). 30 PCR cycle, $1e^{-6}$ PCR mutation rate, PCR efficiency 0.705
603 and 4). HiSeq 2000 100bp sequencing error profile.

604

605 For barcode extraction, the regular expression
606 AAAAAAAAAAGGGGG([ATCG]{14})ATCGATCGTTTTTTT was used in CellBarcode to extract
607 the 14 base-pair random barcode, and the pattern
608 AAAAAAAAAAGGGGGNNNNNNNNNNNNNNNATCGATCGTTTTTTT was used in genBaRcode.
609 Then at the barcode filtering step, the clustering strategy was used, which removed the minority
610 barcodes with a Hamming distance of 1 to the majority ones. We note that CellBarcode discards
611 error reads, whereas genBaRcode adds them to the majority one. We chose this strategy as we
612 found that the resulting underestimation of clone size due to discarding clustered reads was very
613 slight (see comparison of genBaRcode and CellBarcode, Supplementary Figure 1 and 16),
614 whereas if a clustered barcode is actually a real barcode, for example, when library edit distance
615 is small, the result could be a substantial overestimation of some clone sizes. For further
616 information on how this clustering process was carried out, please refer to the methods outlined
617 in the Barcode Filtering section. The run time of above analysis was evaluated by Sys.time
618 function in R 4.2.1⁴⁴. We used CellBarcode version 1.7.1 and genBaRcode version 1.2.6 here
619 and throughout.

620 Acquisition, analysis and simulation of experimental data

621 Several datasets are analyzed in this manuscript; below, for each, we describe first the
622 experimental dataset, then the barcode analysis and finally the simulation parameters (for bulk
623 data).

624 Lentiviral barcode dataset

625 Experimental data

626 We used a lentiviral barcode dataset from our previous publication³⁰. Briefly, it consists of 13,564
627 myeloid cells recovered from mice 4 weeks after transplantation of barcoded EPO-treated
628 HSPCs. The HSPCs were labeled by the LG2.2 barcode library, which has a 20bp fixed-length
629 barcode region, a diversity of > 10000 barcodes, and has a reference library. The myeloid cell
630 DNA was divided into two technical replicates before PCR amplification and sequencing.

631

632 Barcode analysis

633 The output FASTQ file from Eisele et. al.³⁰ was analyzed with the CellBarcode package using the
634 regular expression ACGGAATGCTAGAACACTCGAGATCAG(.{20})ATGTGGTATGATGTATC
635 to extract the 20bp barcode sequence between constant regions. In the regular expression, the

636 first bases ACGGAATG are the plate index used to demultiplex samples with the same P7 index.
637 The extracted barcodes were cleaned by reference library or cluster filtering separately. For the
638 cluster filtering, we remove the minority barcodes with Hamming distance 1 to the majority ones
639 as the barcode library has a minimum edit distance of 5 (Supplementary Figure 11C). Then we
640 normalized the read number (n_i^{reads}) by the total cell count (n_{total}^{cell}) to estimate the clone size
641 (n_i^{cell}) for each barcode clone (i) with following formula:

$$642 \quad n_i^{cell} = n_i^{reads} / \sum_i n_i^{reads} \times n_{total}^{cell} \quad (5)$$

643 For comparing CellBarcode and genBarcode on the fixed-length barcode dataset from Eisele et
644 al. 2022³⁰, both methods use the same criteria to extract and filter barcodes, which involves
645 defining a barcode as a 20bp random sequence between fixed sequences
646 ACGGAATGCTAGAACACTCGAGATCAG and ATGTGGTATGATGTATC. Additionally, cluster
647 filtering is performed to remove minority barcodes with a Hamming distance of 1, the run time was
648 measured by the “Sys.time()” function in R. The spearman correlation was performed using all
649 barcodes.

650

651 Bartender can only define a fixed region of 5bp. Therefore, the barcode definition is set as a 20bp
652 random sequence between ATCAG and ATGTG. The default Bartender clustering filtering has
653 been applied. And the run time was measured by the “time” function in the shell. The Bartender
654 version used here (and following references) is [https://github.com/LaoZZZZZ/bartender-](https://github.com/LaoZZZZZ/bartender-1.1/commit/9683af760cc33f31185140957d503af7f3e230be)
655 [1.1/commit/9683af760cc33f31185140957d503af7f3e230be](https://github.com/LaoZZZZZ/bartender-1.1/commit/9683af760cc33f31185140957d503af7f3e230be).

656

657 Simulation

658 To simulate the barcodes, we used a lentiviral barcode reference library to label 15 cells. The
659 labeled cells were then subjected to clonal expansion, following a log-normal distribution with a
660 mean log clone size of 1.2 and standard deviation of 3. After performing 30 PCR cycles with an
661 error rate of 1e-6, we concatenated the constant regions: 5'
662 ACGGAATGCTAGAACACTCGAGATCAG and 3' ATGTGGTATGATGTATCA. Finally, we
663 simulated the sequencing using the HiSeq2000 profile, aiming for 50 reads per cell.

664 CRISPR gRNA dataset

665 Experimental data

666 Tumor organoids were derived from Apc1638N mice⁵⁰ and transduced with lentiviral particles
667 expressing the Cas9 enzyme along with blasticidin resistance (Addgene plasmid# 52962) as

668 described previously⁵¹. Selection of infected organoids was achieved by adding 10g/ml of
669 blasticidin (#A1113903 Thermofisher) to the medium.

670

671 Cas9-expressing tumor organoids were then transduced with lentiviral particles each containing
672 a sgRNA sequence derived from a bank of 1796 sgRNAs that target Notch1-related genes, as
673 discovered in the study by Mourao et al. 2019⁵². Transduced organoids were harvested either at
674 48h or 7 days post infection. At 7 days, organoids were dissociated and Tomato-expressing live
675 cells (based on DAPI exclusion) were FACS sorted (Supplementary Figure 24). DNA was
676 extracted using a standard phenol:chloroform:isoamyl alcohol protocol. Briefly, cells were
677 resuspended in 500µl of PBS and 1ml of phenol:chloroform:isoamyl alcohol (25:24:1) solution
678 (Sigma #P2069) was added. After centrifugation at 16,000xg for 5 min, the aqueous phase was
679 collected and one volume of chloroform (Sigma #32211) was added. Following a vortex
680 homogenization step, the samples were centrifuged at 16,000xg for 5 min and the aqueous phase
681 was recovered. Precipitation of the DNA was then performed by adding 1µl of glycogen at 20µg/µl
682 (thermofisher # 10814010), 0,5 volume of the sample of 5,5M Sodium Acetate and 2.5 volumes
683 of the sample of cold 100 % Ethanol. After an overnight at -20°C and 30 min of centrifugation at
684 16,000xg at 4°C, the precipitated DNA pellet was recovered in 30µl of water and quantified by
685 nanodrop. 10µl of DNA were then amplified by PCR in triplicates for each sample in order to add
686 P5-staggers and P7-index oligos to perform NGS DNA sequencing. The PCR was performed with
687 Taq polymerase (Promega # M7406) for 22 cycles (30 sec at 95°C, 30 sec at 53°C, 30 sec at
688 72°C).

689

690 The sequences of the primers are the following:

691 P5-staggers:

692 5'AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT[s]T
693 TGTGGAAAGGACGAAACACCG)

694 P7-index:

695 (5'CAAGCAGAAGACGGCATAACGAGATNNNNNNNNGTGACTGGAGTTCAGACGTGTGCTCT
696 TCCGATCTTCTACTATTCTTTCCCTGCACTGT)

697

698 Bead purification of the PCR product using a ratio of 1.2 was performed following the
699 manufacturer's protocol (Beckman Coulter #B23318). Quality and concentration of the samples
700 were assessed on a Tapestation. Then it was sequenced by MiSeq SE110 with 10% PhiX.

701

702 Barcode analysis

703 The gRNA sequencing results are processed by CellBarcode with regular expression
704 “AAGGACGAAACACCG(.{20})”. After reference library-based filtering, the log clone size SD was
705 calculated.

706

707 Simulation

708 We simulated the gRNA sequencing data using a barcode library consisting of 1796 gRNA
709 sequences. The simulated cells were labeled with a clone size log-mean of 1, but varying log
710 clone size SD values ranging from 0.5 to 2.5. To mimic the error rate of Taq polymerase, we
711 performed 20 PCR cycles with a PCR error rate of 10^{-4} . Finally, the sequencing was simulated
712 using the built-in ART MiSeq profile. We analyzed the simulated results in the same manner as
713 the experimental dataset.

714 VDJ barcode MEF cell line dataset

715 Experimental data

716 The VDJ barcodes are produced by an inducible mouse in situ barcode system based on VDJ
717 recombination²⁰. In this system, the V, D and J sequences are separated by the signal cassettes,
718 which are recognized and cut out by the Rag1 (recombination activating gene-1) and Rag2
719 (recombination activating gene-2) enzymes and repaired by non-homologous end joining repair,
720 which is error-prone, creating the diversity of the final barcode sequences. A cassette with
721 reversed *Rag1*, *Rag2* and *TdT* (terminal deoxynucleotidyl transferase) genes are surrounded by
722 LoxP sequences, which can be activated by Cre floxing. The TdT adds de novo nucleotides to
723 the end joins, which increases the diversity of the final barcode sequence.

724

725 In Urbanus et al 2023, mouse embryo fibroblast (MEF) cell lines were created from individual cells
726 of a VDJ barcode-induced mouse with known unique barcode sequences. There are a total of 7
727 MEF cell lines with barcode sequences:

728 CTCGAGGTCATCGAAGTATCAAGTCCAGTTCTACTATCGTAGCTACTA,

729 CTCGAGGTCATCGAAGTATCAAGTCCAGTACTATCGTACTA,

730 CTCGAGGTCATCGAAGTATCAAGTCCAGTCTACTATCGTTACGACAGCTACTA,

731 CTCGAGGTCATCGAAGTATCAAGTCCAGTTCTACTATCGTTACGAGCTACTA,

732 CTCGAGGTCATCGAAGTATCAAGTCCATCGTAGCTACTA,

733 CTCGAGGTCATCGAAGTATCAAGTCCAGTACTGTAGCTACTA,

734 CTCGAGGTCATCGAAGTATCAAGTCCAGTATCGTTACGCTACTA.

735 These cell lines were mixed in specific ratios, in ascending order of powers of 2 from 1 to 7.
736 Sequencing data was then generated with different numbers of initiating cells²⁰.

737

738 Barcode analysis

739 We re-analyzed one of the technical replicates of +0, +2, +4 and +6 PCR cycles with CellBarcode
740 using the regular expression
741 ([ACGT]{12})CTCGAGGTCATCGAAGTATC([ACGT]+)CCGTAGCAAGCTCGAGAGTAGACCTA
742 CT to capture the variable-length barcode between the fixed regions of
743 CTCGAGGTCATCGAAGTATC and CCGTAGCAAGCTCGAGAGTAGACCTACT, after a 12 bp
744 random UMI.

745

746 Simulation

747 For Figure 2B, to simulate a MEF cell line experiment, we simulated 6250 cells (half of the 12,500
748 cells to mimic the technical replicates) with barcode sequences and clone sizes that match the
749 experimental setup. After two cycles of preamplification, a 12bp random UMI is added with a
750 tagging efficiency of 2%. This is followed by 30 cycles of PCR amplification, with a PCR efficiency
751 of 0.705 and a PCR error rate of $1e^{-5}$.

752 For Supplementary Figure 7, we simulated the full dataset to mimic the experiment described
753 above with different numbers of PCR cycles. We used the same barcode sequences, cell number,
754 and type of sequencing while incorporating variable total PCR cycles of +0, +2, +4, and +6. The
755 same fixed 3' sequence as the experimental dataset was added
756 (CCGTAGCAAGCTCGAGAGTAGACCTACTGGAATCAGACCGCCACCATGGTGAGCA
757), and the simulated data were analyzed in the same manner as the experimental dataset.

758 In vivo VDJ barcode mammary gland dataset

759 Experimental data

760 The VDJ barcode mouse was crossed with Notch1Cre^{ERT2} mouse⁵³. Lactating mothers were
761 injected with tamoxifen (0.1mg per g of mouse body mice, MP Biomedicals, 156738) as described
762 ⁵⁴ in order to induce Cre recombination in the progeny at stage P0. Mammary tissue of a DRAG^{+/-}
763 Notch1Cre^{ERT2+/-} female was then collected at 6 weeks of age and mammary single cell
764 dissociation was performed as previously described⁵⁵. Briefly, mammary fat pads were
765 mechanically minced with scissors and scalpel and digested for 90 min at 37C in CO₂-
766 independent medium (Invitrogen, 18045-054) supplemented with 5% fetal bovine serum, 3 mg/ml
767 collagenase A (Roche, 10103586001) and 100 U/ml hyaluronidase (Sigma, H3884). The resulting

768 suspension was sequentially resuspended in 0.25% trypsin–EDTA for 1 min, and then 5 min in 5
769 mg/ml dispase (Roche, 04942078001) with 0.1 mg/ml DNase I (Sigma, D4527) followed by
770 filtration through a 40- μ m mesh. Red blood cells were lysed in NH_4Cl . The obtained single cell
771 suspension was then stained with the following Biolegend antibodies, at a 1/100 dilution: APC
772 anti-mouse CD31 (102510), APC anti-mouse Ter119 (116212), APC anti-mouse CD45 (103112),
773 APC/Cy7 anti-mouse CD49f (313628), and PE anti-mouse EpCAM (118206). Dead cells (DAPI+),
774 and $\text{CD45}^+/\text{CD31}^+/\text{Ter119}^+$ (Lin^+) non-epithelial cells were excluded before analysis using FACS
775 ARIA flow cytometer (Becton Dickinson) (Supplementary Figure 25). In total 20,589 barcoded
776 GFP^+ , Lin^- , $\text{EpCAM}^{\text{high}}$, $\text{CD49f}^{\text{low}}$ luminal cells were sorted into the lysis buffer (Viagen, 301-C).

777

778 For this data we have access to both UMI and non-UMI sequencing libraries as each technical
779 replicate was split in two and processed in parallel with UMI and non-UMI protocols.

780

781 For the UMI barcode sequencing library, we follow the protocol described in²⁰. In brief, the lysed
782 cells were sheared by sonication then divided into two technical replicates, and the target region
783 captured by beads. The DNA in beads was used as a template to do the preamp PCR to amplify
784 the target region with 11 cycles. Next, the UMI was introduced by a second PCR, then the third
785 PCR to add the M1 sequences, finally the fourth PCR to add the adapter sequence to get the
786 sequencing library. The library was sequenced by MiSeq SE110 with 10% PhiX.

787

788 For the non-UMI barcode sequencing library, the preamp PCR product from the UMI barcode
789 library was used to generate a non-UMI sequencing library. We took 100 μ l preamp PCR product,
790 cleaned it with 1.8X SPRI beads, and eluted in 30 μ l DNase free water. The first PCR used 28 μ l
791 of the eluted DNA as template with 50 μ l PCR reaction (10 μ l 5X Q5 buffer, 0.5 μ l 2 U/ μ l Q5 DNA
792 polymerase, 1 μ l 10mM dNTP, 0.25 μ l 100 μ M preamp Fwd primer and preamp Rev primer, 10 μ l
793 DNase free water) for 19 cycles (98C 2min; 19 cycles of 98C 10 sec, 67C 30sec, 72C 30; then
794 72C 5 min). Then the products were cleaned by 1.8 SPRI beads, and eluted into 30 μ l DNase free
795 water. The second PCR used 15 μ l of the eluted DNA from last step as the template with 50 μ l
796 reaction (10 μ l 5X Q5 buffer, 1 μ l 2U/ μ l Q5 DNA polymerase, 1 μ l 10 mM dNTP, 0.25 μ l 100 μ M
797 preamp Fwd primer and M1 Rev primer, and 22.5 μ l DNase free water) for 5 PCR cycles (98C
798 2min; 5 cycles of 98C 10 sec, 67C 30sec, 72C 30; then 72C 5 min). After that, the PCR products
799 were cleaned by 1.8X SPRI beads and eluted into 30 μ l DNase free water. The third PCR used
800 10 μ l DNA from last step as template to add the NGS adaptors by 20 μ l PCR reaction (4 μ l 5X Q5
801 buffer, 0.4 μ l 2U/ μ l Q5 DNA polymerase, 0.4 μ l 10mM dNTP, 0.1 μ l 100 μ M P5 tagging primer, 4 μ l

802 2.5uM P7 tagging primer with index, and 1.1 ul DNase free water) by 5 PCR cycles (98C 2min; 5
803 cycles of 98C 10 sec, 67C 30sec, 72C 30; then 72C 5 min). The final DNA was cleaned by 1X
804 SPRI beads, and eluted into 30ul DNase free water. The library with 10% PhiX was sequenced
805 by MiSeq in SE110 mode with a 25M sequencing chip aimed for 20M reads output. This library
806 was sequenced together with other samples but independent to the UMI barcode library.

807

808 preamp Fwd primer:

809 ACTCACTATAGGGAGACGCGTGTTACC

810 preamp Rev primer:

811 GACACGCTGAACTTGTGGCCGTTTA

812 M1 Rev primer:

813 AGTTCAGACGTGTGCTCTTCCGATCCAGCTCGACCAGGATGGG

814 P5 tagging primer:

815 AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCTACT

816 CACTATAGGGAGACGCGTGTT

817 P7 tagging primer:

818 CAAGCAGAAGACGGCATAACGAGATTGACTGAGTGACTGGAGTTCAGACGTGTGCTCTTCC

819 GATC

820

821 Barcode analysis

822 For the VDJ barcode UMI library analysis using CellBarcode, we extract the barcode and UMI
823 using regular expression
824 “(.{16})CCTCGAGGTCATCGAAGTATCAAG(.*?)CCGTAGCAAGCTCGAGAGTAGACCTACT”,
825 which defines the 16bp UMI sequence before the constant region and the variable-length VDJ
826 barcode sequence followed by another constant region. Then we removed the UMI-barcode tags
827 with fewer than 100 reads and counted the UMIs per barcode in the remaining tags, which is a
828 robust threshold since the final barcode is very stable when we increase the threshold
829 (Supplementary Figure 18A), and used the remaining barcodes. Investigators can use a similar
830 approach to determine a read count threshold, in conjunction with knowledge about their targeted
831 sequencing depth per cell.

832

833 For the Non-UMI barcode library sequencing, we used the regular expression
834 “CCTCGAGGTCATCGAAGTATCAAG(.*?)CCGTAGCAAGCTCGAGAGTAGACCTACT” to match
835 the variable-length VDJ barcode between the constant regions. The automatic read count
836 threshold was used to identify true barcodes.

837

838 We compared CellBarcode and Bartender using both UMI and non-UMI sequencing described
839 above. We extracted the barcodes with CellBarcode as described above, with the UMI tag
840 requiring a minimum of 100 reads to be counted. In non-UMI libraries, an automatic threshold is
841 applied in CellBarcode. For Bartender, it only allows a maximum of a 5bp match in the fixed
842 region. Therefore, the barcode is defined between the fixed regions TCAAG and CCGTA. The
843 UMI is defined by the first 16bp random sequence in both cases. Then the clustering with 1
844 mismatch is used for both UMI and non-UMI sequencing. The run time of Bartender is measured
845 by shell command “time”, and for CellBarcode by the “Sys.time()” function in R. The shared
846 barcodes were counted and visualized using a Venn plot. Linear regression was performed on
847 the shared barcodes.

848

849 Simulation

850 We used CellBarcodeSim to simulate the above VDJ sequencing data. The simulation included
851 a VDJ barcode library with 100 cells, which were expanded using a log-normal distribution (log
852 clone size mean 1.2, SD 1). We used a random UMI of length 16bp, and sequenced 100 reads
853 per UMI using the ART built-in MiSeq profile, resulting in sequences of length 111 bp. Additionally,
854 we added fixed regions at the 5' end (CCTCGAGGTCATCGAAGTATCAAG) and the 3' end
855 (CCGTAGCAAGCTCGAGAGTAGACCTACTGGAATCAGACCGCCACCATGGTGAGCACACG
856 TCTGAACTCCAGTCACTCAGTCAATCTCGTATGCCGTCTTCTGCTTG). Other parameters
857 were kept default.

858 CellTag barcode scRNA-seq dataset

859 Experimental data

860 The scRNA-seq CellTag BAM file (Bidy et. al. 2018)³⁶ was downloaded from SRA with access
861 number SRR7347033. This file corresponds to the Mouse Embryonic Fibroblast (MEF) cell line
862 that was infected with CellTag barcodes, underwent fate reprogramming through overexpression
863 of transcription factors FOXA1 and HNF4 α , and was sequenced after 15 days.

864

865 Barcode analysis

866 For the CellTagR analysis, we followed its demo described here [https://github.com/morris-](https://github.com/morris-lab/CellTagR)
867 [lab/CellTagR](https://github.com/morris-lab/CellTagR). Firstly, we filtered the BAM file in bash by 1). Filtering unmapped reads, and 2).
868 Filtering transgene reads. The filtered BAM file was used as input to both the CellTagR and
869 CellBarcode pipelines. After first creating a CellTag object, the V1 barcode was extracted from

870 the BAM file, by matching 5' constant GGT and 3' constant GAATTC. After that, barcode filtering
871 was applied including: 1). Filter cells (a list of cells passing QC was downloaded from GEO with
872 dataset id GSE99915) 2). Barcode sequence error correction with clustering using Starcode, 3).
873 Keep UMIs with at least 2 reads, and 4). Barcode reference library filtering (whitelist filtering). The
874 barcode reference library (whitelist) can be found with the demo datasets of the CellTagR
875 package. Barcode clustering error correction was done by starcode-1.4⁵⁶.

876

877 We applied the CellTagR pipeline described above as closely as possible using CellBarcode.
878 Using CellBarcode, we extracted the V1 barcode using the regular expression
879 "GGT([ATCG]{8})GAATTC" which matches the 8bp DNA sequence surrounded by two fixed
880 constant regions. Then, we carried out the 4 filtering steps using the CellBarcode package which
881 are 1). Filter cells using the QC passed list described above, 2). Barcode sequencing correction
882 by removing minority barcodes with a Hamming distance of 1 to the majority one, 3). Keep UMI
883 with at least 2 reads and 4). Barcode reference library filtering.

884 VDJ barcode scRNA-Seq dataset

885 *Barcode analysis*

886 In this section we describe VDJ barcode extraction with CellBarcode, the barcode filtering was
887 described in the Results section.

888 In single cell sequencing data analysis, each cell is stored as an individual sample in the
889 BarcodeObj, and this object has the same data structure as that of bulk analysis.

890 The FASTQ file was acquired from the authors. Their read 1 and read 2 were concatenated. In
891 the sequence, we defined the cellular 10X barcode as the first 16 bases, and the UMI as 12 bases
892 followed, according to the 10X 3' scRNA-seq reads structure. And the lineage barcode sequence
893 was extracted using the 3' and 5' constant sequences: "CGAAGTATCAAG" and
894 "CCGTAGCAAG".

895

896 The result in original paper was accessed from GitHub: [https://github.com/TeamPerie/Cosgrove-](https://github.com/TeamPerie/Cosgrove-et-al-2022/blob/main/Figure1/RNA_BC_PREPROCESSING/input_output_m534/agrep_10xbc_and_v)
897 [et-al-](https://github.com/TeamPerie/Cosgrove-et-al-2022/blob/main/Figure1/RNA_BC_PREPROCESSING/input_output_m534/agrep_10xbc_and_v)

898 [2022/blob/main/Figure1/RNA_BC_PREPROCESSING/input_output_m534/agrep_10xbc_and_v](https://github.com/TeamPerie/Cosgrove-et-al-2022/blob/main/Figure1/RNA_BC_PREPROCESSING/input_output_m534/agrep_10xbc_and_v)
899 [bc_m534_both.txt.gz](https://github.com/TeamPerie/Cosgrove-et-al-2022/blob/main/Figure1/RNA_BC_PREPROCESSING/input_output_m534/agrep_10xbc_and_v). A brief description of the barcode filtering of the original is as follows: UMIs
900 were filtered to keep only those with 3 or more reads and one dominant VDJ barcode (defined
901 as ≥ 0.45 reads). The dominant barcode for each UMI was extracted, and finally they assigned
902 one VDJ barcode to a 10x cell if there is good agreement across UMIs, defined as ≥ 0.75

903 agreement across all remaining UMIs. If there is only one UMI retained, they further ensured that
904 the VDJ barcode for this UMI was the dominant barcode across all the reads for that cell and
905 has ≥ 0.45 of reads.

906

907 **Statistics and Reproducibility**

908 No statistical method was used to predetermine sample size. No data were excluded from the
909 analyses. The experiments were not randomized. The Investigators were not blinded to allocation
910 during experiments and outcome assessment.

911 Data availability

912 The lentiviral barcodes dataset from Eisele et al (2022)³⁰ was obtained from ⁵⁷; the corresponding
913 pre-analysed data is available at: <https://github.com/TeamPerie/Eisele-et-al>. The CellTag
914 barcode sequencing data from Bidy et al (2018)³⁶ is on GEO with dataset ID GSE99915. The
915 Marsolier et al (2022)¹⁷ barcoded scRNASeq dataset is on GEO with dataset ID GSE164716. The
916 mammary gland VDJ barcode dataset and gRNA sequencing data is available on Zenodo⁵⁸. The
917 MEF cell line mixes VDJ barcode dataset is available from ⁵⁹. The VDJ-barcoded scRNA-seq data
918 from Cosgrove et al (2023) belongs to the authors of that paper and was given to us for the
919 purposes of this paper; to obtain this data please contact Leïla Perié (leila.perie@curie.fr). Source
920 data for Figures 2-6 are provided with this paper.

921 Code availability

922 Code for all analysis in this study is available at
923 https://github.com/TeamPerie/CellBarcode_paper_Sun_et_al and at ⁶⁰. The CellBarcode
924 package is available at Bioconductor
925 <https://bioconductor.org/packages/release/bioc/html/CellBarcode.html> and at ⁶¹.
926 <https://doi.org/doi:10.18129/B9.bioc.CellBarcode>). And the Barcode sequencing simulation kit is
927 available at <https://github.com/TeamPerie/CellBarcodeSim> and at ⁶².

928 Acknowledgements

929 We would like to acknowledge the valuable discussions with the members of the Perié and Fre
930 labs and the NGS facility. We appreciate the assistance from the animal facility, flow cytometry
931 and NGS facilities at Institut Curie. We would like to express our gratitude to Emily Tubeuf and
932 Cecile Conrad from the Perié lab for their assistance in conducting experiments and analyzing
933 FACS data, respectively. The study was supported by the European Research Council (ERC)
934 under the European Union's Horizon 2020 research and innovation programme ERC StG 758170-
935 Microbar (to L.P.), an ATIP Avenir grant from CNRS and Bettencourt-Schueller Foundation (to
936 L.P.) and by PSL* Research University (to S.F.), the French National Research Agency (ANR)
937 (ANR-15-CE13-0013-01) (to S.F.), FRM Equipes (EQU201903007821) (to S.F.), FSER
938 (Schlumberger Foundation) (FSER20200211117) (to S.F.), ARC Foundation label 2022 (n°
939 ARCPGA2021120004232_4874) (to S.F.) and Labex DEEP (ANR-11-LBX-0044) (to S.F.).

940 Author Contributions Statement

941 Conceptualization was carried out by W.S., A.M.L., and L.P., while W.S., M.P., M.H., S.F., L.P.
942 and M.F. were responsible for barcode experiment design, protocol development and performed
943 the experiments. The analysis and programming were completed by W.S., who also managed
944 data curation and generated the figures. Writing was a collaborative effort of W.S., L.P. and
945 A.M.L.. L.P. and A.M.L. provided supervision throughout the study, and S.F. and M.F. contributed
946 to the review and editing process.

947

948 Competing Interests Statement

949 The authors declare no competing interests.

950

951 Figure legends

952 **Figure 1.** CellBarcode: a package to extract and identify lineage barcodes.

953 A. Barcode experiment scheme. Cells are labeled with genetic barcodes, divide and differentiate,
954 with progeny inheriting the barcode. Barcodes are read out by next generation sequencing (NGS)
955 in descendant cells. CellBarcode allows extraction, filtering and identification of barcodes from
956 NGS data and returns a barcode count matrix for further analysis.

957 B. Diagram of barcode sequencing data processing with CellBarcode. CellBarcode reads the raw
958 sequencing data (FASTQ, FASTA, BAM/SAM files or R object) and checks the QC (QC and
959 filtering functions) before extracting the barcode sequences (barcode extraction functions).
960 Barcodes are then filtered to remove PCR and sequencing errors using different filtering
961 strategies (barcode cleaning functions). After filtering, barcode data can be plotted with the visual
962 check functions and exported as a barcode frequency matrix (export functions).

963 C. Example of barcode processing workflow using CellBarcode. Barcodes (underlined) are
964 extracted from raw sequences using a regular expression (sequence in bold) that depends on the
965 barcode type. Barcodes are then filtered, as detailed in D, to eliminate spurious barcodes and
966 exported.

967 D. The four most commonly used barcode filtering strategies. Green indicates true barcodes, red
968 spurious barcodes. 1) Reference library filtering: barcodes B1, B2 and B3 that match the
969 reference list are considered true barcodes, M3 and M5 are removed. 2) Threshold filtering:
970 barcodes that have a read number superior or equal to the threshold of 20 are kept (B1 and B2)

971 and barcodes below the threshold are removed (M3, M5 and B3). 3) cluster filtering: barcodes
972 with an edit distance smaller than a threshold to a more abundant barcode are eliminated. Here,
973 two barcodes have one substitution difference (mutant loci in white) from an abundant barcode
974 and will be deleted. 4) UMI filtering: usually involves retaining the most abundant sequence per
975 UMI followed by a UMI count threshold per barcode.

976 **Figure 2.** Cellular barcode sequencing simulation

977 A. Schematic of barcoding experiment simulation with CellBarcodeSim and the parameters that
978 can be tuned at each step, starting with simulation of a barcode library, cell labeling and clonal
979 expansion, PCR amplification and finally sequencing. The round shape represents
980 undifferentiated cells, the triangle and rectangle represent differentiated cell types.

981 B, C. Stacked bar plots, created using CellBarcode, displaying the percentage of bases for the
982 VDJ barcode dataset with UMI (B) and a random barcode dataset (C) across each sequencing
983 cycle. Each column represents a sequencing cycle, with color and height indicating the base and
984 proportion respectively. Both simulated and real experimental data are presented for each
985 dataset. The percentage of total reads matching the regular expression is indicated, as well as
986 the Pearson correlation between the most abundant base per sequencing cycle. Fixed and/or
987 UMI regions are annotated above the heatmap. The VDJ barcode dataset is the MEF line
988 experiment data with 12500 cells from Urbanus et al. 2022²⁰; the random barcode dataset is from
989 Eisele et al. 2022³⁰. Simulation details for each dataset are provided in the Methods section.

990 **Figure 3.** Benchmarking Barcode Filtering Strategies with Simulated Data.

991
992 A, B. Percentage precision and recall of true barcodes for different threshold filtering using read
993 proportion thresholds of 0.0001 (A), and 0.001 (B). Several scenarios with two types of barcodes
994 (random and VDJ) and three different clone size variations across barcodes are compared.

995 C. Area Under Precision-Recall Curve (PR-AUC) using threshold filtering for two types of
996 barcodes (random and VDJ) and three different clone size variations across barcodes.

997 D. Same as in C after reference filtering.

998 E. Same as in C after cluster filtering.

999 F. Diagrams depicting reference library filtering and cluster filtering advantages and drawbacks.

1000 Reference library filtering removes spurious barcodes that are not in the library but keeps spurious
1001 barcodes that match a barcode in the reference library. Cluster filtering removes low abundance

1002 barcodes that are similar to abundant barcodes. This can result in the removal of true barcodes
1003 which have sequence similarity to another true barcode, for example if the barcode library has
1004 small edit distance.

1005 G. PR-AUC after UMI filtering for variable-length VDJ barcodes for two higher clone size variations
1006 (log clone size standard deviation of 2 and 3). An initial filtering based on UMI count greater than
1007 10 reads was performed before computing PR-AUC.

1008 H. Barcode filtering decision tree.

1009 Except when otherwise specified, each simulated scenario has the reference parameters from
1010 Supplementary Table 2: 30 simulations, 300 induced barcodes with log clone size mean 1.2, PCR
1011 cycle 30, PCR efficiency 0.705, PCR error 1×10^{-6} , reads per cell 50 and sequencing profile
1012 HiSeq2000. Specifically for H, the number of PCR cycles before and after UMI tagging are 10 and
1013 20 respectively, with 8 bp UMI and tagging efficiency 0.02. The median and interquartile range
1014 (IQR, the difference between the 75th and 25th percentiles of the data) are shown in the boxplot
1015 over 30 simulations, and the outliers (beyond the whiskers of $Q3 + 1.5IQR$ or $Q2 - 1.5IQR$) plotted
1016 as dots.

1017 The two-sided Wilcoxon Test is applied to compare the precision, recall or AUC of different
1018 simulation conditions.

1019 **Figure 4.** Lentiviral barcode sequencing analysis.

1020

1021 A. Base quality heatmap made with CellBarcode. Each row is a sample, each column corresponds
1022 to a sequencing cycle, the color represents the median base Phred quality score.

1023 B. Base percentage plotted against the sequencing cycle number made with CellBarcode.
1024 Sequence shows a 20bp barcode with fixed flanking regions either side. Color represents a base
1025 pair.

1026 C. Barcode normalized read count + 1 (by total 10^5 read) as filtered in the original paper Eisele
1027 et al. (2022)³⁰ versus using CellBarcode. Each dot is a barcode. The spearman correlation and
1028 p-value (two-sided) are displayed in the top left corner.

1029 D. Barcode cell counts between the two technical replicates for the data without filtering. The read
1030 counts were normalized to cell counts. Each dot is a barcode with black indicating presence in
1031 the reference library provided in Eisele et al. (2022)³⁰.

1032 E. Same as D but after cluster filtering, the filtering process involves removing barcodes that have
1033 a Hamming distance of less than 2 from a more abundant barcode.

1034 In C, D, and E the red line represents $y = x$, the black line indicates a threshold of one cell.
1035

1036 **Figure 5.** In vitro VDJ barcode analysis.

1037 A. Sequencing library design and sequencing scheme. A sample was divided into two technical
1038 replicates. After a first PCR amplification, each technical replicate was further divided into two for
1039 sequencing library preparation with and without UMIs.

1040 B. Stacked bar plot made with CellBarcode showing the base percentage for each sequencing
1041 cycle. Each column corresponds to a sequencing cycle, the color and height indicate the base
1042 and proportion respectively. Both rows depict the same biological sample, with or without UMI for
1043 sequencing. The position of the regular expression (constant region) and the UMI are annotated.

1044 C. Barcode read counts between technical replicates for the No-UMI library without filtering.
1045 Automatic thresholds (marked by red lines) were applied to remove the errors in each technical
1046 replicate separately. The numbers show the barcode count in each of the four categories as
1047 divided by the threshold lines. Each dot represents a barcode. Plot made with CellBarcode, the
1048 dots are semi-transparent to display overlap.

1049 D. Barcode UMI count between technical replicates with UMI library. The data was first filtered
1050 retaining UMI with at least 10 reads. The red lines indicate a UMI count threshold of 1. The number
1051 of barcodes in each of the four categories as divided by the threshold lines is annotated. Each
1052 dot represents a barcode. Plot made with CellBarcode, the dots are semi-transparent to display
1053 overlap.

1054 E. Comparing the number of barcodes identified in the No-UMI library and the UMI library in one
1055 technical replicate. For the No-UMI library, the automatic threshold was applied as shown in C.
1056 For the UMI library, the same filtering steps were applied as in D with the addition of a UMI count
1057 threshold of 1.

1058 F. Barcode read count after filtering between the No-UMI library and the UMI library for one of the
1059 technical replicates. The read counts were renormalized to one. A linear regression was fitted,
1060 and the fitted line (and shaded area of 95% confidence interval) and its parameters are written on
1061 the plot. Each dot represents a barcode.

1062 **Figure 6.** Single cell RNASeq Cellular DNA barcode analysis.

1063 A. Diagram of how single cell sequencing lineage barcode sequencing data are processed with
1064 CellBarcode. Input files can be FASTQ or BAM/SAM files. The lineage barcodes are extracted,
1065 filtered, and exported for subsequent analysis.

1066 B. Filtering steps for single cell sequencing lineage barcode data implemented in CellBarcode.
1067 Firstly, for each UMI the dominant barcode is identified and other barcodes are removed; then
1068 UMIs with a read count below a threshold are removed. For each barcode, the number of UMI is
1069 counted and the barcodes are filtered based on a UMI count threshold.

1070 C. The number of lineage barcodes found per cell before and after filtering barcodes based on
1071 the dominant barcode per UMI using the VDJ scRNA-seq data from Cosgrove et al³⁵so. In the
1072 scatter plot the y-axis is the barcode number in a cell, each dot repr³⁵esents a cell, and the
1073 distribution is shown by the violin plot.

1074 D. The number of lineage barcodes per cell (corresponding to the left black y-axis) and the cell
1075 number (corresponding to the right red y-axis) for different thresholds of read per UMI. The data
1076 was first processed with the dominant barcode per UMI filter. Each black dot represents a cell
1077 while the violin plot shows the distribution of the barcode number per cell.

1078 E. The number of lineage barcodes per cell (corresponding to the left black y-axis) and the number
1079 of cells with a unique barcode (corresponding to the right red y-axis) for different thresholds of
1080 UMI count per barcode. The data was first processed with the dominant barcode per UMI filter
1081 and the UMI read threshold ≥ 2 . In the figure, each dot represents a cell, and the distribution is
1082 shown by the violin plot. The red line plot represents the number of retained, unique barcoded
1083 cells after applying different UMI count filters described in x-axis.

1084 **References**

- 1085 1. Sankaran, V. G., Weissman, J. S. & Zon, L. I. Cellular barcoding to decipher clonal dynamics
1086 in disease. *Science* **378**, eabm5874 (2022).
- 1087 2. Perié, L. & Duffy, K. R. Retracing the in vivo haematopoietic tree using single-cell methods.
1088 *FEBS Lett.* **590**, 4068–4083 (2016).
- 1089 3. Lu, R., Neff, N. F., Quake, S. R. & Weissman, I. L. Tracking single hematopoietic stem cells
1090 in vivo using high-throughput sequencing in conjunction with viral genetic barcoding. *Nat.*
1091 *Biotechnol.* **29**, 928–933 (2011).
- 1092 4. Kok, L., Masopust, D. & Schumacher, T. N. The precursors of CD8+ tissue resident memory
1093 T cells: from lymphoid organs to infected tissues. *Nat. Rev. Immunol.* **22**, 283–293 (2022).
- 1094 5. Naik, S. H. *et al.* Diverse and heritable lineage imprinting of early haematopoietic progenitors.
1095 *Nature* **496**, 229–232 (2013).
- 1096 6. Dhimolea, E. *et al.* An Embryonic Diapause-like Adaptation with Suppressed Myc Activity
1097 Enables Tumor Treatment Persistence. *Cancer Cell* **39**, 240-256.e11 (2021).
- 1098 7. Merino, D. *et al.* Barcoding reveals complex clonal behavior in patient-derived xenografts of
1099 metastatic triple negative breast cancer. *Nat. Commun.* **10**, 766 (2019).
- 1100 8. Echeverria, G. V. *et al.* Resistance to neoadjuvant chemotherapy in triple negative breast
1101 cancer mediated by a reversible drug-tolerant state. *Sci. Transl. Med.* **11**, eaav0936 (2019).
- 1102 9. Echeverria, G. V. *et al.* High-resolution clonal mapping of multi-organ metastasis in triple
1103 negative breast cancer. *Nat. Commun.* **9**, 5079 (2018).
- 1104 10. Blundell, J. R. & Levy, S. F. Beyond genome sequencing: Lineage tracking with barcodes to
1105 study the dynamics of evolution, infection, and cancer. *Genomics* **104**, 417–430 (2014).
- 1106 11. Naik, S. H., Schumacher, T. N. & Perié, L. Cellular barcoding: a technical appraisal. *Exp.*
1107 *Hematol.* **42**, 598–608 (2014).

- 1108 12. McKenna, A. *et al.* Whole-organism lineage tracing by combinatorial and cumulative genome
1109 editing. *Science* **353**, (2016).
- 1110 13. Frieda, K. L. *et al.* Synthetic recording and in situ readout of lineage information in single cells.
1111 *Nature* **541**, 107–111 (2017).
- 1112 14. Alemany, A., Florescu, M., Baron, C. S., Peterson-Maduro, J. & van Oudenaarden, A. Whole-
1113 organism clone tracing using single-cell sequencing. *Nature* **556**, 108–112 (2018).
- 1114 15. Raj, B., Gagnon, J. A. & Schier, A. F. Large-scale reconstruction of cell lineages using single-
1115 cell readout of transcriptomes and CRISPR–Cas9 barcodes by scGESTALT. *Nat. Protoc.* **13**,
1116 2685–2713 (2018).
- 1117 16. Spanjaard, B. *et al.* Simultaneous lineage tracing and cell-type identification using CRISPR–
1118 Cas9-induced genetic scars. *Nat. Biotechnol.* **36**, 469–473 (2018).
- 1119 17. Marsolier, J. *et al.* H3K27me3 conditions chemotolerance in triple-negative breast cancer.
1120 *Nat. Genet.* **54**, 459–468 (2022).
- 1121 18. Thielecke, L. *et al.* Limitations and challenges of genetic barcode quantification. *Sci. Rep.* **7**,
1122 1–14 (2017).
- 1123 19. Pei, W. *et al.* Polylox barcoding reveals haematopoietic stem cell fates realized in vivo. *Nature*
1124 **548**, 456–460 (2017).
- 1125 20. Urbanus, J. *et al.* DRAG in situ barcoding reveals an increased number of HSPCs contributing
1126 to myelopoiesis with age. *Nat. Commun.* **14**, 2184 (2023).
- 1127 21. Beltman, J. B. *et al.* Reproducibility of Illumina platform deep sequencing errors allows
1128 accurate determination of DNA barcodes in cells. *BMC Bioinformatics* **17**, 151 (2016).
- 1129 22. Lyne, A.-M. *et al.* A track of the clones: new developments in cellular barcoding. *Exp. Hematol.*
1130 **68**, 15–20 (2018).
- 1131 23. Hadj Abed, L., Tak, T., Cosgrove, J. & Perié, L. CellDestiny: A RShiny application for the
1132 visualization and analysis of single-cell lineage tracing data. *Front. Med.* **9**, (2022).

- 1133 24. Espinoza, D. A., Mortlock, R. D., Koelle, S. J., Wu, C. & Dunbar, C. E. Interrogation of clonal
1134 tracking data using barcodetrackR. *Nat. Comput. Sci.* **1**, 280–289 (2021).
- 1135 25. Lin, D. S. *et al.* DiSNE Movie Visualization and Assessment of Clonal Kinetics Reveal Multiple
1136 Trajectories of Dendritic Cell Development. *Cell Rep.* **22**, 2557–2566 (2018).
- 1137 26. Thielecke, L., Cornils, K. & Glauche, I. genBaRcode: a comprehensive R-package for genetic
1138 barcode analysis. *Bioinformatics* **36**, 2189–2194 (2020).
- 1139 27. Zhao, L., Liu, Z., Levy, S. F. & Wu, S. Bartender: a fast and accurate clustering algorithm to
1140 count barcode reads. *Bioinformatics* **34**, 739–747 (2018).
- 1141 28. Kong, W. *et al.* CellTagging: combinatorial indexing to simultaneously map lineage and
1142 identity at single-cell resolution. *Nat. Protoc.* **15**, 750–772 (2020).
- 1143 29. Bandler, R. C. *et al.* Single-cell delineation of lineage and genetic identity in the mouse brain.
1144 *Nature* **601**, 404–409 (2022).
- 1145 30. Eisele, A. S. *et al.* Erythropoietin directly remodels the clonal composition of murine
1146 hematopoietic multipotent progenitor cells. *eLife* **11**, e66922 (2022).
- 1147 31. Sender, R. & Milo, R. The distribution of cellular turnover in the human body. *Nat. Med.* **27**,
1148 45–48 (2021).
- 1149 32. Bystrykh, L. V. Generalized DNA Barcode Design Based on Hamming Codes. *PLOS ONE* **7**,
1150 e36852 (2012).
- 1151 33. Beneyto-Calabuig, S. *et al.* Clonally resolved single-cell multi-omics identifies routes of
1152 cellular differentiation in acute myeloid leukemia. *Cell Stem Cell* **30**, 706-721.e8 (2023).
- 1153 34. Jindal, K. *et al.* Multiomic single-cell lineage tracing to dissect fate-specific gene regulatory
1154 programs. 2022.10.23.512790 Preprint at <https://doi.org/10.1101/2022.10.23.512790> (2022).
- 1155 35. Cosgrove, J. *et al.* Metabolically Primed Multipotent Hematopoietic Progenitors Fuel Innate
1156 Immunity. 2023.01.24.525166 Preprint at <https://doi.org/10.1101/2023.01.24.525166> (2023).
- 1157 36. Bidy, B. A. *et al.* Single-cell mapping of lineage and identity in direct reprogramming. *Nature*
1158 **564**, 219–224 (2018).

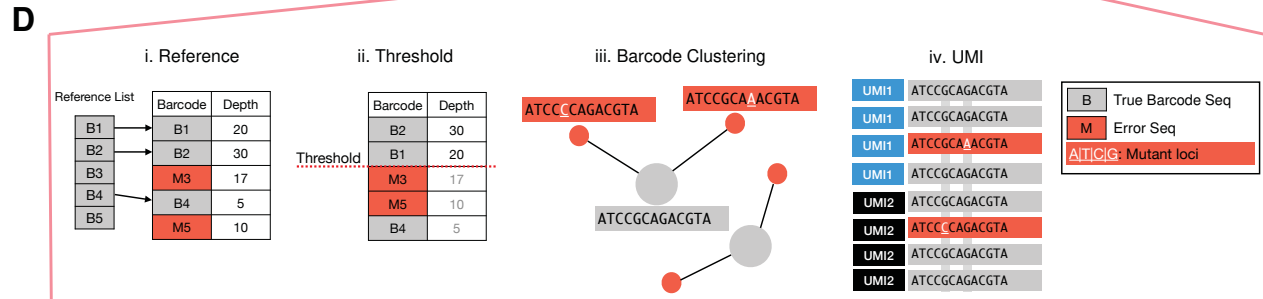
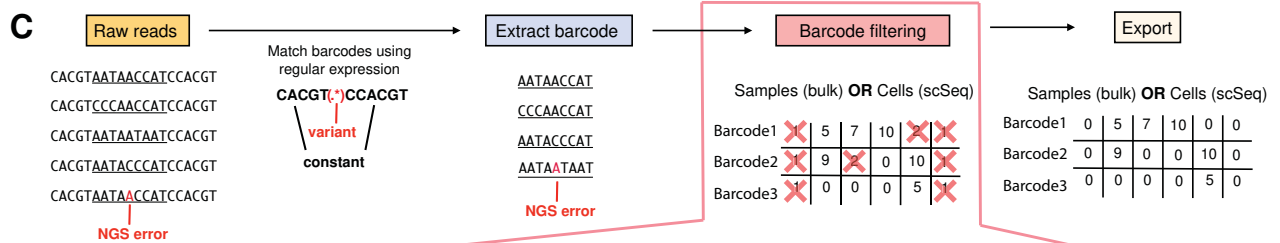
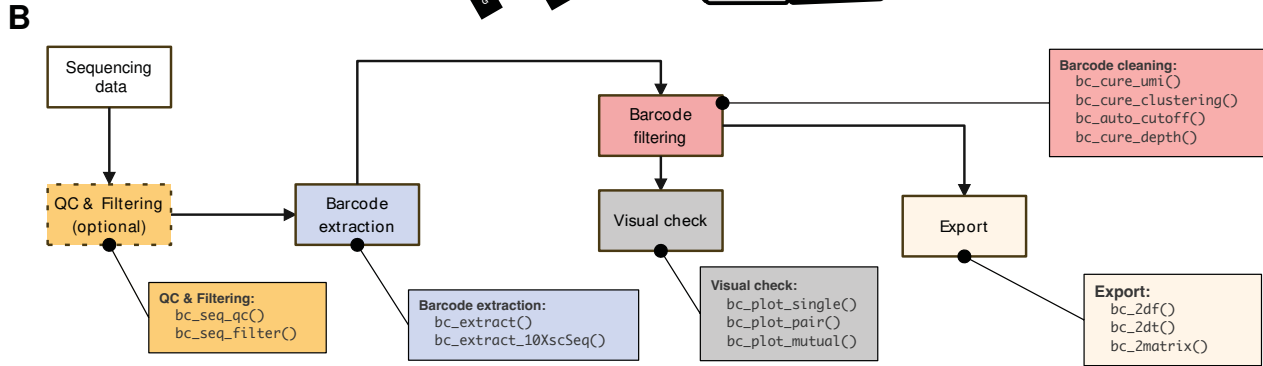
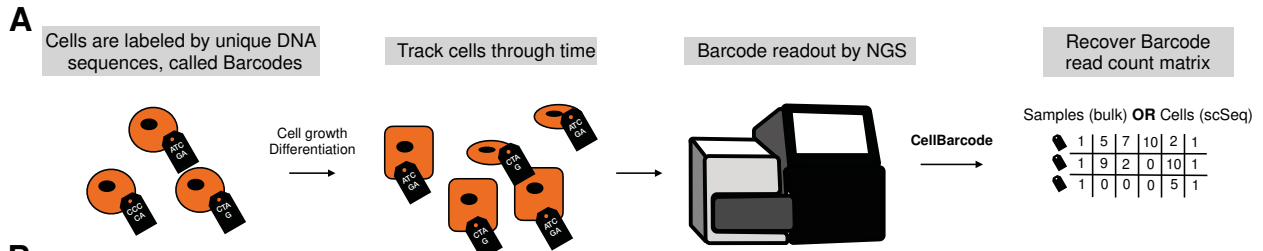
- 1159 37. Radtke, S. *et al.* Stochastic fate decisions of HSCs after transplantation: early contribution,
1160 symmetric expansion, and pool formation. *Blood* **142**, 33–43 (2023).
- 1161 38. Error Detecting and Error Correcting Codes - Hamming - 1950 - Bell System Technical
1162 Journal - Wiley Online Library. [https://onlinelibrary.wiley.com/doi/10.1002/j.1538-](https://onlinelibrary.wiley.com/doi/10.1002/j.1538-7305.1950.tb00463.x)
1163 [7305.1950.tb00463.x](https://onlinelibrary.wiley.com/doi/10.1002/j.1538-7305.1950.tb00463.x).
- 1164 39. Huang, W., Li, L., Myers, J. R. & Marth, G. T. ART: a next-generation sequencing read
1165 simulator. *Bioinformatics* **28**, 593–594 (2012).
- 1166 40. Buschmann, T. DNABarcodes: an R package for the systematic construction of DNA sample
1167 tags. *Bioinforma. Oxf. Engl.* **33**, 920–922 (2017).
- 1168 41. Marcou, Q., Mora, T. & Walczak, A. M. High-throughput immune repertoire analysis with
1169 IGoR. *Nat. Commun.* **9**, 1–10 (2018).
- 1170 42. Desponds, J., Mora, T. & Walczak, A. M. Fluctuating fitness shapes the clone-size distribution
1171 of immune repertoires. *Proc. Natl. Acad. Sci.* **113**, 274–279 (2016).
- 1172 43. Adair, J. E. *et al.* DNA Barcoding in Nonhuman Primates Reveals Important Limitations in
1173 Retrovirus Integration Site Analysis. *Mol. Ther. - Methods Clin. Dev.* **17**, 796–809 (2020).
- 1174 44. Team, R. C. R: A language and environment for statistical computing. R Foundation for
1175 Statistical Computing, Vienna, Austria. *HttpwwwR-Proj.* (2016).
- 1176 45. Weiss, G. & von Haeseler, A. A coalescent approach to the polymerase chain reaction.
1177 *Nucleic Acids Res.* **25**, 3082–3087 (1997).
- 1178 46. McInerney, P., Adams, P. & Hadi, M. Z. Error Rate Comparison during Polymerase Chain
1179 Reaction by DNA Polymerase. *Mol. Biol. Int.* **2014**, 287430 (2014).
- 1180 47. Sing, T., Sander, O., Beerenwinkel, N. & Lengauer, T. ROCR: visualizing classifier
1181 performance in R. *Bioinformatics* **21**, 3940–3941 (2005).
- 1182 48. Wang, H. & Song, M. Ckmeans.1d.dp: Optimal k-means Clustering in One Dimension by
1183 Dynamic Programming. *R J.* **3**, 29–33 (2011).

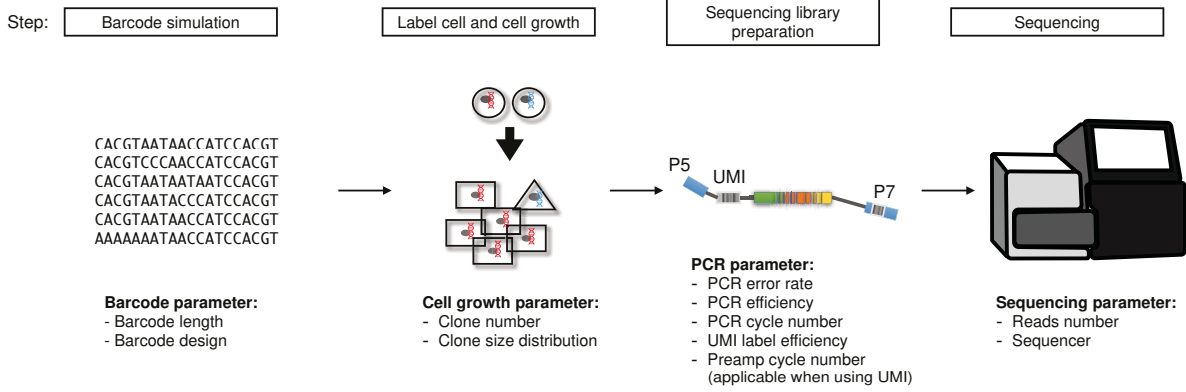
- 1184 49. Johnson, M. S., Venkataram, S. & Kryazhimskiy, S. Best Practices in Designing, Sequencing,
1185 and Identifying Random DNA Barcodes. *J. Mol. Evol.* **91**, 263–280 (2023).
- 1186 50. Fodde, R. *et al.* A targeted chain-termination mutation in the mouse *Apc* gene results in
1187 multiple intestinal tumors. *Proc. Natl. Acad. Sci.* **91**, 8969–8973 (1994).
- 1188 51. Jacquemin, G. *et al.* Paracrine signalling between intestinal epithelial and tumour cells
1189 induces a regenerative programme. *eLife* <https://elifesciences.org/articles/76541/figures>
1190 (2022) doi:10.7554/eLife.76541.
- 1191 52. Mourao, L. *et al.* Lineage tracing of Notch1-expressing cells in intestinal tumours reveals a
1192 distinct population of cancer stem cells. *Sci. Rep.* **9**, 888 (2019).
- 1193 53. Fre, S. *et al.* Notch Lineages and Activity in Intestinal Stem Cells Determined by a New Set
1194 of Knock-In Mice. *PLOS ONE* **6**, e25785 (2011).
- 1195 54. Lilja, A. M. *et al.* Clonal analysis of Notch1-expressing cells reveals the existence of unipotent
1196 stem cells that retain long-term plasticity in the embryonic mammary gland. *Nat. Cell Biol.* **20**,
1197 677–687 (2018).
- 1198 55. Lloyd-Lewis, B. *et al.* In vivo imaging of mammary epithelial cell dynamics in response to
1199 lineage-biased Wnt/ β -catenin activation. *Cell Rep.* **38**, 110461 (2022).
- 1200 56. Zorita, E., Cuscó, P. & Fillion, G. J. Starcode: sequence clustering based on all-pairs search.
1201 *Bioinformatics* **31**, 1913–1919 (2015).
- 1202 57. Eisele, A. S. *et al.* Erythropoietin directly remodels the clonal composition of murine
1203 hematopoietic multipotent progenitor cells. (2021) doi:10.5281/zenodo.5645045.
- 1204 58. SUN, W. *et al.* CellBarcode package paper dataset. (2023) doi:10.5281/zenodo.8124948.
- 1205 59. Urbanus, J. *et al.* UrbanusCosgrove-et-al-DRAG-mouse. (2023)
1206 doi:10.5281/zenodo.10027001.
- 1207 60. Sun, W. *et al.* TeamPerie/CellBarcode_paper_Sun_et_al. (2024)
1208 doi:10.5281/zenodo.10492761.
- 1209 61. Sun, W. *et al.* CellBarcode. doi:10.18129/B9.bioc.CellBarcode.

1210 62. Sun, W. *et al.* TeamPerie/CellBarcodeSim. (2024) doi:10.5281/zenodo.10492831.

1211

1212



A**B**