

1 **CloneRetriever: retrieval of rare clones from heterogeneous cell populations**

2

3 David Feldman<sup>1,2†</sup>, FuNien Tsai<sup>1†</sup>, Anthony J. Garrity<sup>1</sup>, Ryan O'Rourke<sup>1,3</sup>, Lisa Brenan<sup>1</sup>,  
4 Patricia Ho<sup>1,3</sup>, Elizabeth Gonzalez<sup>1,3</sup>, Silvana Konermann<sup>4</sup>, Cory M. Johannessen<sup>1\*</sup>,  
5 Rameen Beroukhim<sup>1,6,7\*</sup>, Pratiti Bandopadhyay<sup>1,5,6\*</sup>, Paul C. Blainey<sup>1,8\*</sup>

6

7 1. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA.

8 2. Department of Physics, MIT, Cambridge, Massachusetts 02142, USA.

9 3. Dana-Farber/Boston Children's Cancer and Blood Disorders Center, Boston,

10 Massachusetts 02115, USA.

11 4. Salk Institute, La Jolla, California 92037, USA.

12 5. Department of Pediatrics, Harvard Medical School, Boston, Massachusetts 02115,

13 USA.

14 6. Division of Medical Oncology, Dana-Farber Cancer Institute, Boston, Massachusetts

15 02115, USA.

16 7. Department of Medicine, Harvard Medical School, Boston, Massachusetts 02115,

17 USA.

18 8. Department of Biological Engineering, MIT, Cambridge, Massachusetts 02142, USA

19

20 †These authors contributed equally to this work

21 \* The correspondence should be addressed to: [pblainey@broadinstitute.org](mailto:pblainey@broadinstitute.org),

22 [pratiti\\_bandopadhyay@dfci.harvard.edu](mailto:pratiti_bandopadhyay@dfci.harvard.edu), [rameen@broadinstitute.org](mailto:rameen@broadinstitute.org), and

23 [cory.johannessen@gmail.com](mailto:cory.johannessen@gmail.com)

24 A.J.G., R.O., and C.M.J. have current addresses at Arbor Biotechnologies, Casma

25 Therapeutics, and Novartis Institutes for BioMedical Research, respectively.

## 26 **Abstract**

## 27 **Background**

28 Many biological processes, such as cancer metastasis, organismal development, and  
29 development of resistance to cytotoxic therapy, rely on the emergence of rare sub-  
30 clones from a larger population. Understanding how the genetic and epigenetic features  
31 of diverse clones affect clonal fitness provides insight into molecular mechanisms  
32 underlying selective processes. However, identifying causal drivers of clonal fitness  
33 remains challenging. Population-level analysis has limited resolution to characterize  
34 clones prior to selection, while high-resolution single-cell methods are destructive and  
35 challenging to scale across large populations, limiting further functional analysis of  
36 relevant clones.

## 37 **Results**

38 Here, we develop CloneRetriever, a methodology for tracking and retrieving rare clones  
39 throughout their response to selection. CloneRetriever utilizes a CRISPR sgRNA-  
40 barcode library that allows isolation of viable cells from specific clones within the  
41 barcoded population using a sequence-specific retrieval reporter. We demonstrated that  
42 CloneRetriever can measure clonal fitness of cancer cell models *in vitro* and retrieve  
43 targeted clones at abundance as low as 1 in 1,883 in a heterogeneous cell population.

## 44 **Conclusions**

45 CloneRetriever provides a means to track and access specific and rare clones of  
46 interest across dynamic changes in population structure to comprehensively explore the  
47 basis of these changes.

48

## 49 **Keywords**

50 Cellular heterogeneity, Viable clone-specific cells recovery, Clonal fitness tracking,  
51 CRISPR sgRNA-barcode DNA library

52

## 53 **Introduction**

54 The response of a heterogeneous population to selection pressure is shaped by the  
55 growth dynamics of individual clones within the population. Rare clones can play a  
56 decisive role in the outcome of selection. Examples include evasion of anti-retroviral  
57 therapy by rare HIV variants [1], expansion of drug-resistant cancer cells under  
58 chemotherapy [2], and seeding of metastases by clonal tumor cells [3], [4]. Studying  
59 how genetic and epigenetic differences affect the fitness of different clones during  
60 selection provides an opportunity to understand both how the selective process  
61 operates and how populations are reshaped by selection. In particular, identifying

62 causal drivers of clone fitness could give rich insights into the molecular mechanisms of  
63 selection and suggest potential interventions.

64 Heritable and plastic cellular features can drive selection outcomes. For example,  
65 genetic features can change with mutagens, such as DNA-damaging chemotherapies,  
66 and epigenetic states can rapidly shift in response to drug exposure [5] or environment  
67 [6]. Metastatic clones may alter their epigenetic profiles upon seeding a metastatic site  
68 [7], obscuring the preexisting features that enabled them to metastasize. However,  
69 existing methods to identify these features tend to rely on comparing populations in bulk  
70 before and after selection, which limits their usefulness in detecting pre-existing features  
71 that changed during selection. A useful alternative approach would be to identify clones  
72 based upon their response to selective pressure, and then isolate representative  
73 untreated cells from each clone for genomic and functional characterization.

74 Genomically integrated DNA barcodes provide a scalable methodology to track rare  
75 clones by measuring relative barcode abundance over time [8]. However, relative clone  
76 fitness alone cannot elucidate mechanisms of selection. Single-cell technologies can  
77 provide genomic profiles of heterogeneous cells within a population. Clone identity can  
78 be incorporated into single-cell RNA-seq (scRNA-seq) profiles by capturing transcribed  
79 barcodes, linking clonal history and cell fate [9]. However, single-cell genomic profiling  
80 is inherently destructive. Both DNA barcoding and single-cell approaches have a limited  
81 ability to probe functional differences between clones, whereas retrieval of viable cells  
82 from clones would enable a wide range of genomic and functional analysis.

83 Here, we report CloneRetriever, an experimental system that permits tracking,  
84 selection, and recovery of arbitrarily chosen, viable clones from a cell population.  
85 CloneRetriever employs a diverse library of single-guide RNAs (sgRNAs). In the  
86 absence of Cas9 activity, these serve as inert barcodes for tracking cells. In the  
87 presence of Cas9, these sgRNAs direct Cas9 in a clone-specific fashion to activate a  
88 reporter. Cas9-dependent reporter expression permits the physical isolation of specific  
89 cells within a population while preserving cell viability. This methodology allows for the  
90 isolation and comparative analysis of specific clones at any stage of evolution. Isolated  
91 cells can then be characterized by downstream functional assays, such as phenotypic  
92 characterization, genetic perturbation, or small molecule screens, thus enabling  
93 comprehensive analysis of how clonal features affect fitness.

94

## 95 **Results**

### 96 **Overview of the barcoding and retrieval strategy**

97 To enable tracking and retrieval of clones within a heterogeneous population, we  
98 designed a selectable barcode strategy that allows for retrieval of viable cells with  
99 clone-specific barcodes. In this system, each clone is tagged with a library of random  
100 CRISPR sgRNAs [10]. In the absence of Cas9 expression, the sgRNA-barcodes serve  
101 as inert labels that are propagated upon cell division, similar to previously reported  
102 clonal barcoding strategies [8],[5]. The relative abundance of each clone can be  
103 quantified by deep sequencing of the DNA-integrated sgRNA-barcode. The relative

104 fitness of clones can then be determined by sequencing sgRNA-barcodes over time  
105 (e.g., before and after drug selection). By expanding the ancestral barcoded population  
106 and splitting the daughter cells into replicate selection assays, clone-specific fitness  
107 differences can be estimated (e.g., clones with a drug-dependent fitness advantage)  
108 (Figure 1a).

109 We designed this system so that specific clones can be isolated from a barcoded  
110 population using a retrieval vector with a target site matching the sgRNA-barcode of  
111 interest (Figure 1b). Introducing Cas9 nuclease leads to double-strand DNA breaks at  
112 the target site specifically in the clone that expresses the corresponding sgRNA-  
113 barcode. DNA repair generates frameshift mutations at the target site, which may shift  
114 the translation frame of one or more downstream reporters [11] (Figure 1c). Activation of  
115 the retrieval reporter can result in both gain and loss of reporter expression (e.g., a shift  
116 that brings a GFP reporter into frame and an RFP reporter out of frame).

117

## 118 **An sgRNA-barcode library enables tracking clonal subpopulations**

119 We generated two high complexity sgRNA-barcode libraries using fully degenerate  
120 oligonucleotide templates of either 20- or 26- nucleotides (nt) (Additional file 1: Fig.  
121 S1a). To test the clone tracking capacity of sgRNA-barcodes, we applied the 26-nt  
122 barcode library to monitor clonal resistance to the BET-bromodomain inhibitor JQ1, in  
123 D458, a MYC-amplified medulloblastoma cell line known to contain pre-existing  
124 resistant clones to a chemotherapeutic [5]. We first transduced D458 cells with the 26-nt

125 barcode library at low MOI (< 0.3). We then selected the transduced cells with  
126 puromycin and restricted the population size to ensure that a high fraction of barcodes  
127 corresponded to unique clones (Figure 2a).

128 We expanded the barcoded D458 population and split it into replicates that were treated  
129 with either 2  $\mu$ M JQ1 or DMSO only (vehicle control). Deep sequencing at the time of  
130 the replicate split (early time point, or ETP) detected 84,014 barcodes prior to drug  
131 selection (Figure 2a). After 52 days, we harvested cells and quantified barcode  
132 abundance in each replicate (Additional file 3: Table S3). An average of 2,938 barcodes  
133 were enriched in JQ1-treated replicates, comprising about 2% of the original barcodes.  
134 Approximately 50% of the JQ1-selective resistant barcodes were shared by all  
135 replicates (Figure 2b and c), suggesting both that these barcodes marked clones with  
136 predetermined resistance to JQ1 and that our barcode library enables tracking of clones  
137 with such heritable phenotypes within a heterogeneous population. Analysis of barcode  
138 enrichments showed no significant biases based on barcode GC content or homology to  
139 the genome, suggesting that the sgRNA-barcode library can function similarly to other  
140 inert barcoding libraries (Additional file 1: Fig. S1b and c). Although the 26nt-barcoding  
141 library enables tracking complex clonal populations, previous studies showed that  
142 increasing the length of gRNA sequences above 20nt leads to reduced Cas9 activity  
143 [12]. Therefore, we opted to employ a 20-nt sgRNA barcode library that we have shown  
144 is also able to track evolution of populations through targeted therapies [13].

145

## 146 **Design of a retrieval vector activated by frameshift mutations**

147 To retrieve viable clones, we designed a frameshift reporter that can be specifically  
148 activated by an sgRNA-barcode of interest. This approach relies on the generation of  
149 insertion or deletion (indel) mutations by Cas9 nuclease in a target region to shift the  
150 translation frame of a reporter cassette, similar to vectors used to monitor gene-editing  
151 outcomes [11], [14]. An alternative approach would be to use CRISPR-a (dCas9-  
152 transcriptional activator) to activate marker expression in a barcode-dependent fashion  
153 [15]. However, we found that a transcriptional activation-based reporter lacked  
154 specificity, in part due to a high background level of transcription in a fraction of cells  
155 subsequent to lentiviral integration of the reporter (Additional file 1: Fig. S2a and b).  
156 Conversely, frameshift reporters have the potential for extremely high specificity due to  
157 the low background rate of activating mutations. We opted to deliver the reporter using  
158 a lentiviral system, as it can effectively transduce a wide range of cell lines. Lentiviral  
159 transduction at low MOI followed by antibiotic selection integrates a single reporter copy  
160 into most cells, minimizing the potential for a cell to contain multiple reporters in different  
161 frameshift states.

162 We designed a retrieval vector that gains GFP fluorescence in response to a +2  
163 frameshift mutation that occurs within a narrow targeting window of ~100 bp. The vector  
164 contains two cassettes respectively in the +0 and +2 translation frames: a selection  
165 marker (e.g., blasticidin, Blast) and a fluorescent protein (e.g., mCherry) linked by a T2A  
166 self-cleaving peptide in the +0 frame, and a second fluorescent protein (e.g., GFP) in  
167 the +2 frame. The +0 cassette (mCherry-T2A-Blast) is located downstream of the +2



168 cassette in order to aid in selecting for integrants with the correct initial frame via  
169 antibiotic selection (blasticidin) or fluorescence-activated cell sorting (FACS) (mCherry)  
170 (Figure 1c). To minimize the likelihood of background activation, we included triple stop  
171 codons in all reading frames immediately upstream of the Kozak translation initiation  
172 site. All sequences downstream of the translation initiation site were codon optimized to  
173 eliminate start and stop codons that could interfere with reporter performance  
174 (Methods).

175 In order to target a specific barcode, the matching target sequence is cloned into the  
176 targeting window between the translation start site and the beginning of the GFP coding  
177 sequence. Targeting of Cas9 nuclease by the sgRNA-barcode generates indel  
178 mutations in the targeting window. If a +2 indel occurs, the reading frame shifts such  
179 that the +0 cassette is out of frame, while the +2 cassette is in frame, giving rise to GFP  
180 expression (Figure 1c). In addition to GFP, a variety of alternative selection elements,  
181 such as antibiotic resistance or surface affinity markers, can be used to assist in  
182 enriching for cells with +2 frameshifts.

183

#### 184 **The retrieval vector is specifically activated by target sgRNA-barcodes**

185 We further applied two modifications to improve the activity of our retrieval vector. With  
186 the initial version (TMv1), activation with the matching guide produced >1% GFP  
187 compared to 0.001% mismatch guide controls (Additional file 1: Fig. S3). To improve  
188 sensitivity, we replaced GFP with mNeonGreen and switched the EFS promoter to a

189 stronger EF1a promoter (TMv2). To allow FACS-independent enrichment, we also  
190 expanded the +0 selection cassette to include either a Zeocin resistance or H2K surface  
191 affinity marker upstream of mNeonGreen (TMv2-Zeo, TMv2-H2K). Compared to TMv1,  
192 the TMv2 retrieval vector showed approximately 10-fold increased sensitivity at  
193 comparable specificity.

194 We then systematically evaluated the performance of TMv2 using 5 randomly selected  
195 barcodes from our sgRNA-barcode library and matching targets cloned into TMv2 and  
196 TMv2-Zeo. We generated HeLa-TetR-Cas9 cell lines expressing each individual  
197 sgRNA-barcode, so that specificity and sensitivity could be directly assessed by flow  
198 cytometry (Figure 3a). All five barcodes activated mNeonGreen expression from a  
199 matching retrieval vector (Figure 3b). In the same experiment, each retrieval vector was  
200 tested with mismatched barcode targets to evaluate specificity (Figure 3a). The results  
201 showed a low false positive rate ranging from 0 to  $2.7 \cdot 10^{-5}$  for TMv2, and from 0 to  
202  $5.5 \cdot 10^{-4}$  for TMv2-Zeo. The sensitivity for the matched barcodes ranged from  $1.8 \cdot 10^{-1}$   
203 to  $2.3 \cdot 10^{-2}$  for TMv2, and from  $0.92 \cdot 10^{-1}$  to  $3.8 \cdot 10^{-2}$  for TMv2-Zeo, suggesting that the  
204 system was capable of high specificity and selectivity.

205 In addition to single barcode reporters, multiplexed activation of several barcodes with  
206 one reporter can be achieved by expanding the target sequence to contain targets for  
207 multiple sgRNA-barcodes (Figure 3a). To demonstrate multiplexing, we designed  
208 retrieval vectors to target three independent sgRNA-barcode sequences (Figure 3c).  
209 These vectors showed similar sensitivity to those individual sgRNA-barcodes, albeit at

210 2.6-fold reduced specificity ( $1.4 \cdot 10^{-3}$ ), possibly due to the increased likelihood of  
211 background mutations in the expanded target region.

212

### 213 **Identification and viable isolation of rare hygromycin-resistant HeLa cells**

214 We next tested our ability to retrieve drug-resistant and drug-sensitive clones of interest  
215 in a well-controlled setting. We engineered hygromycin-resistant HeLa-TetR-Cas9 cells  
216 and spiked them into a pool of hygromycin-sensitive HeLa-TetR-Cas9 cells to achieve a  
217 final population of cells in which 2% of all cells expressed the hygromycin resistance  
218 gene.

219 We transduced the cells with the 20-nt sgRNA-barcode library at low MOI, and then  
220 bottlenecked, expanded, and cryopreserved them in replicate vials. Sequencing of one  
221 replicate verified the presence of 441 barcodes ranging in abundance from 1 in 100 to 1  
222 in 100,000 (Figure 4a, Barcoding). To assay for hygromycin resistance, we split the  
223 cells and treated them in replicate with either hygromycin or PBS (vehicle) (Figure 4a,  
224 Selection). We then nominated candidate hygromycin-resistant barcodes by comparing  
225 the abundance of the barcodes in hygromycin-treated cells to the PBS-treated groups  
226 (Figure 4a, Deconvolution). We found 15 candidate hygromycin-resistant barcodes with  
227 >10-fold enrichment after hygromycin treatment and an ETP frequency of at least 1 in  
228 3,000.

229 We carried out retrieval for 4 clonal barcodes: one hygromycin-sensitive barcode  
230 candidate (T1) and 3 hygromycin-resistant barcode candidates (T2, T3 and T4) that  
231 were represented in the population at frequencies ranging from 1 in 652 (T2) to 1 in  
232 140,000 (T4) (Figure 4b and d) (Additional file 1: Table S1 and Additional file 4: Table  
233 S2). We also analyzed 2 types of control populations: cells transduced with retrieval  
234 vectors targeting barcodes not present in the library, and cells without doxycycline  
235 induction of Cas9.

236 For each sgRNA-barcode, we cloned a matching targeting sequence into the retrieval  
237 vectors TMv2 and TMv2-Zeo. To retrieve cells representative of the initial, unselected  
238 population, we thawed and expanded barcoded cells preserved at the ETP. Barcoded  
239 cells were transduced with either TMv2 or TMv2-Zeo and selected with blasticidin for 4  
240 days. Blasticidin was then removed and Cas9 expression was induced with doxycycline  
241 for 7 days (Figure 4c, Retrieval vector transduction).

242 This process greatly enriched clones T1-T3. FACS purification followed by expansion  
243 and sequencing indicated up to 845-fold enrichment of these clones relative to the ETP  
244 fraction, to a minimum purity of 44.87% (T3) and a maximum purity of 92.51% (T1)  
245 (Figure 4d and e, FACS). In addition to the FACS-based enrichment, we also carried  
246 out selection using TMv2-Zeo, and detected 12 - 85-fold enrichment (Figure 4e, Zeocin  
247 selection). Clone T4 was not present in the enriched population, suggesting the  
248 sensitivity of the retrieval vectors was insufficient to recover viable clones present at  
249 frequencies in the population that are smaller than 1 in 140,000 (Figure 4e).

250 **Validation of retrieved clones and analysis of sensitivity-limiting background**  
251 **events**

252 In order to confirm that the barcoding and retrieval protocols led to the recovery of  
253 clones exhibiting the hygromycin-resistant phenotype, we sorted individual cells  
254 transduced with TMv2 into multi-well plates and expanded them as clonal populations  
255 (Figure 4c, Clone enrichment & isolation). We analyzed a total of 132 single-cell clones  
256 by deep sequencing of their sgRNA-barcodes. We detected 2 populations: clones with  
257 exact matches to the targeted barcodes (52/132 clones), and mismatched clones  
258 (barcode edit distance >9; 80/132 clones) (Figure 4f and Additional file 1: Fig. S4a, red  
259 square). The hygromycin sensitivity of the single-cell clones reflected their barcodes,  
260 with the exception of one single-cell clone with a candidate hygromycin resistance  
261 barcode (T3) that was sensitive to hygromycin (Figure 4g and Additional file 1: Fig. S4a,  
262 blue square).

263 To investigate activation of the retrieval vector, we performed Sanger sequencing of 75  
264 clones over a 2-kb region encompassing the translation start site, barcode-specific  
265 targeting region, and mNeonGreen coding sequence. As expected, clones with the  
266 correct sgRNA-barcode contained +2 frameshift mutations in the targeting region, with a  
267 distribution of indel sizes consistent with repair by non-homologous end-joining following  
268 Cas9 cleavage (Additional file 1: Fig. S4c) [16]. In contrast, 21/43 false positive clones  
269 exhibited a ~80-nt stereotyped deletion immediately upstream of the mNeonGreen  
270 coding sequence (Additional file 1: Fig. S4b and c). Deletions due to lentiviral intra-  
271 molecular recombination between homologous regions are well-characterized [17].

272 However, our codon-optimized retrieval vector lacks substantial homology near the  
273 deleted region (no repeated kmers with length >7), suggesting an alternative  
274 mechanism. The false positive events observed were largely due to the stereotyped  
275 deletion, as we found that sorting error likely did not contribute to these false positive  
276 events, as re-analysis of expanded clones showed that all clones contained  
277 GFP+/mCherry- cells (Additional file 1: Fig. S4a, green square).

278 Together, these results indicate that CloneRetriever is capable of tracking hygromycin-  
279 resistant phenotypes under treatment, and enriching rare clones up to 800-fold.

280

## 281 **Discussion**

282 We engineered a molecular tool that couples an sgRNA-barcode library for tracking  
283 clones with a Cas9-based frameshift reporter to isolate viable cells representing target  
284 clones from the population. A challenge to studying the mechanisms underlying clonal  
285 evolution has been that bulk methods have limited resolution to observe characteristics  
286 of rare clones, while single-cell methods have not been able to target clones with known  
287 evolutionary paths. We showed that our system can accurately track clonal fitness  
288 under drug selection and allows efficient retrieval of a targeted set of clones at  
289 frequencies as low as 1 in 1,883. We demonstrate that a CRISPR sgRNA-barcode  
290 approach is able to scale to high complexity libraries capable of barcoding  $>10^5$  clones,  
291 while a frameshift retrieval reporter activated by barcode-specific Cas9-mediated  
292 mutations enables fluorescence-based retrieval of clones. The sgRNA-barcode design

293 is especially conducive to multiplexing for simultaneous retrieval of a handful of clones  
294 at a time (as shown in Figure 3c), because it allows straightforward expansion of the  
295 activating window to accommodate multiple sgRNA-barcodes.

296 Isolating clonally barcoded cells from an untreated, ancestral population enables direct  
297 testing of mechanisms underlying differential clone fitness. Unlike bulk methods that rely  
298 on strong positive selection to enrich for cells of interest, our method allows retrieval  
299 from clones with any fitness profile, such as slow-growing, persistent, or negatively  
300 selected clones. The ability to expand pure populations of target clones enables the use  
301 of a broad range of functional and molecular profiling assays. For example, access to  
302 pure populations enables high-input assays to determine how epigenetic alterations,  
303 such as changes in DNA methylation and chromatin state, affect fitness differences  
304 between genetically similar clones. Deep characterization of purified resistant clones  
305 can be useful in identifying resistant drivers, and through perturbational approaches, the  
306 association between these putative drivers and phenotype can be defined. The sgRNA-  
307 barcodes can also be readily adapted to existing high-throughput single-cell readouts  
308 developed for CRISPR screens, such as single-cell gene expression [18], [19] and  
309 optical screening [20].

310 An important caveat for this and other lentiviral-based DNA barcoding strategies is the  
311 possibility of unintended side effects from semi-random lentiviral integration on  
312 barcoded clones. Lentiviral integration can either disrupt or increase gene activity,  
313 leading to clone-specific effects. Our approach, in which clones of interest are isolated,  
314 simplifies sequencing of the DNA barcode insertion site, which can help rule out

315 integration-driven effects. While introducing the retrieval vector requires an additional  
316 lentiviral integration event, multiple independent sub-clones can be retrieved per clone  
317 of interest, serving as biological replicates for the retrieval process. A second potential  
318 confounding factor is sgRNA-barcode sequence-specific effects on clone fitness in the  
319 absence of Cas9. However, we found no significant correlation between enrichment in a  
320 clone tracking experiment and sgRNA-barcode sequence content or sgRNA homology  
321 to the genome.

322 Background activation of the frameshift retrieval reporter may hinder applications where  
323 the clones of interest exist at low frequency. In principle, a reporter activated by an indel  
324 mutation in a 100-nt activating window could have a background rate as low as ~1 in 1  
325 billion per cell division, as the rate of naturally occurring indel mutations in human cells  
326 is estimated to be ~1 in  $10^{11}$  indel/bp/cell division per generation [21]-[23]. We identified  
327 stereotyped deletions in the T2A linker region as the primary source of false positive  
328 activations; optimizing this sequence could significantly suppress background.

329 Alternatively, negative selection against the GFP-containing frame could be applied  
330 prior to editing to remove cells with a premature frameshift. To improve sensitivity to  
331 both +1 and +2 frameshifts, a second reporter cassette could be added in the +1 frame.  
332 Selecting sgRNA-barcodes based on their predicted indel distribution could further  
333 increase activation efficiency [24].

334



## 335 **Conclusions**

336 Clone tracking and retrieval enable deep, mechanistic studies in a wide range of  
337 selection scenarios. For example, tracking cells during reprogramming or differentiation  
338 protocols would enable isolation and epigenetic characterization of ancestral clones that  
339 are predisposed to successful outcomes [9], [25]. Similarly, retrieving untreated cells  
340 from clones surviving mutagenic chemotherapy, such as alkylating agents, could  
341 address outstanding questions about whether resistance is pre-existing or acquired [26].  
342 Clones can also be targeted based on fitness profiles derived from multiple parallel  
343 selection conditions. Altogether, live clone retrieval capability enables barcoding  
344 experiments to advance from observing clone frequency statistics toward  
345 experimentally-driven mechanistic studies by providing access to key samples  
346 supporting a wide range of genomic and functional assays.

347

## 348 **Methods**

### 349 **Library construction**

350 Degenerate oligos for sgRNA-barcode library construction were synthesized by IDT and  
351 cloned into lentiGuide-Puro [27] by Gibson assembly as previously reported [28].  
352 Approximately 300 ug of Gibson product was transformed into 25 uL of Endura  
353 electrocompetent cells (Lucigen). After a 1 hour recovery period, 0.1% of transformed  
354 bacteria were plated in a 10-fold dilution series on ampicillin plates to determine the

355 number of successful transformants. The remainder of the transformed bacteria were  
356 cultured in 50 mL of LB with 50 ug/mL ampicillin for 16 hours at 30° C. Plasmid libraries  
357 were extracted using Plasmid MidiPlus kit (Qiagen) and sequenced to a depth of 95  
358 million reads on Illumina Nextseq, corresponding to 13X coverage of 3.9 million  
359 barcodes. Lentivirus was prepared as previously reported [28] by transfecting a total of  
360 10 million HEK 293FT cells. The library virus was determined by transduction and  
361 puromycin selection in HeLa-Tet-Cas9 cells to contain 600 million infective particles,  
362 corresponding to a 153X coverage of barcodes.

363

#### 364 **Barcoding of HeLa and D458 cell lines**

365 HeLa-Tet-Cas9 cells were cultured in DMEM medium supplemented with 10%  
366 tetracycline-screened FBS (Hyclone) and 1% penicillin-streptomycin. sgRNA-barcodes  
367 were transduced as previously described [28] and selected with 1 ug/mL puromycin for  
368 3 days. The lentiviral multiplicity of infection (MOI) was determined to be between 0.05  
369 and 0.3 for all libraries, so that a majority of cells carry a single integrated sgRNA-  
370 barcode. Barcoded cell lines were expanded to a total of  $1.0 \cdot 10^7$  cells and  
371 cryopreserved in aliquots of  $1.0 \cdot 10^6$  cells for subsequent drug selection and retrieval.  
372 D458 medulloblastoma cells were cultured in DMEM/F12 media supplemented with  
373 10% FCS and 1% GPS (glutamate, pen-strep). Four million cells were transduced with  
374 the sgRNA barcode library (10 wells of  $3.0 \cdot 10^6$  cells with 50ul of virus) by spin infection

375 (1,000g, 120 minutes, 30° C). Selection with 1 ug/mL puromycin was initiated 48 hours  
376 post-transduction and maintained for a total of 3 days.

377

### 378 **Drug resistance experiments: D458 and JQ1**

379 Barcoded D458 medulloblastoma cells (fingerprint verified) were treated with DMSO or  
380 JQ1 at a concentration of 2uM in multiple replicate plates (5 x DMSO and 5 x JQ1).  
381 Four million barcoded D458 cells were plated in each replicate plate in presence of  
382 DMSO or JQ1. Barcoded D458 cells were also frozen in 10% DMSO/FCS for future  
383 retrieval. In addition, cells were collected for DNA-extraction to determine barcode  
384 representation at the early-time point (ETP). Cells were retreated with compound every  
385 3-4 days. Cells were counted and passaged every 3-4 days, maintaining a minimum  
386 representation of 4 million cells. Cells were cultured in DMSO or JQ1 for a total of 52  
387 days prior to harvesting for DNA extraction for barcode sequencing and deconvolution.

388

### 389 **Drug resistance experiments: HeLa and hygromycin**

390 HeLa-TetR-Cas9 cells were infected with a lentiviral ORF construct  
391 (pLX\_TRC317\_PGK-Hygro) containing a hygromycin resistance cassette. After  
392 selection with 300 ug/ml hygromycin for 1 week, HeLa-LacZ cells were spiked into  
393 uninfected cells at a ratio of 1:50. Cells were then infected with the CloneRetriever  
394 library at MOI <0.3. Following selection with puromycin, we plated a fixed number of

395 cells (to achieve a ‘bottleneck’ of the number of barcoded cells) and expanded the  
396 population. Cells were frozen in liquid nitrogen (early time point, ETP) in replicates of  
397  $1 \cdot 10^7$  cells. One replicate was thawed for barcoding experiments (1 x ETP, 5 x DMSO  
398 and 5 x hygromycin at 300 ug/ml). Replicate cells were cultured in DMSO or hygromycin  
399 for 16 days, after which DNA was extracted from both the ETP control and  
400 DMSO/hygromycin treated replicates for barcode sequencing and deconvolution. At  
401 each passage, we ensured the number of cells plated was at least 10-fold the library  
402 complexity in order to maintain representation.

403

#### 404 **Library deconvolution**

405 Genomic DNA was extracted and prepared for deep sequencing as reported [28].  
406 Libraries were sequenced to a minimum depth of 18 million reads, corresponding to a  
407 barcode coverage of >80X. Counts of sgRNA-barcodes were obtained by filtering for  
408 reads containing exact matches to the flanking sequences, and matches with <3 reads  
409 were discarded.

410

#### 411 **Clonal fitness measurements**

412 Relative clone abundances were calculated from normalized read counts and clones  
413 were ranked by abundance within each replicate. For the D458 clonal tracking  
414 experiment, JQ1- and DMSO-enriched barcodes were defined as those with a median

415 rank above 4,000 in JQ1 replicates and above 2,000 in DMSO replicates. NGS data  
416 analysis were run with Python 2.7 with its libraries numpy 1.13.1, matplotlib 2.1.2,  
417 seaborn 0.9.9, and jupyter 4.3.0.

418

### 419 **Retrieval reporter construct**

420 The mNeon, T2A, Zeocin, H2K, and Blastidicin coding sequences were codon  
421 optimized with silent nucleotide substitutions to remove out-of-frame start and stop  
422 codons. Oligos containing targeting barcode sequences and PAM (NGG) matching  
423 barcodes of interest were synthesized (IDT) and cloned into frameshift reporter  
424 plasmids by golden gate assembly. All targeting barcode sequences were filtered to  
425 have <70% GC content, no more than 4 consecutive repeated bases and no stop  
426 codons. Lentivirus was prepared as previously described [28] and transduced into  
427 barcoded HeLa-Tet-Cas9 cells at an MOI of <0.3. After 4 days of selection with 10  
428 ug/mL blastidicin, 1 ug/mL doxycycline was added to induce Cas9 expression. Cells  
429 were harvested for deep sequencing as previously reported [28].

430

### 431 **FACS sample preparation and analysis**

432 HeLa cells were carefully washed with PBS and trypsinized with TrypLE Express  
433 (Gibco) for 5 minutes. DMEM media contained 10% FBS and 1%

434 Penicillin/Streptomycin was used to neutralize trypsin prior to FACS analysis.  
435 Fluorescent protein expression was measured on a Cytoflex flow cytometer. FlowJo  
436 V10 was used for analysis. Populations were sorted with high-purity mode on a SONY-  
437 SH800 FACS machine, and expanded for 2 weeks before deep sequencing. All  
438 analyzed populations were first gated on FSC-A/FSC-H and FSC-A/SSC-A to identify  
439 singlets and cells respectively (Additional file 1: Fig. S5).

440

#### 441 **Characterization of clones**

442 FACS-sorted clones were trypsinized in plate 3 days after sorting and further expanded  
443 for ~7 days. GFP or mCherry expression for each clone was validated on a Cytoflex  
444 flow cytometer. To determine hygromycin sensitivity, the clones were treated with or  
445 without 300 ug/ml hygromycin and the media was replenished with fresh hygromycin  
446 every 3 days for 7 days. Cell number was measured with a Cytoflex flow cytometer. For  
447 Sanger analysis, a 2-kb region of the lentiviral transgene was PCR-amplified from the  
448 EF1a promoter (forward strand, primer pTM\_negative\_fwd) and from the Blast gene  
449 (reverse strand, primer pTM\_negative\_rev) and sequenced with sanger sequencing  
450 primer (pTM\_sanger\_primer) (Additional file 1: Table S2).

451

## 452 **Analysis of frameshift status and indel calculation**

453 For each clone, we used the corresponding unedited retrieval vector as a reference  
454 sequence for alignment of Sanger sequencing traces. We determined the location of  
455 insertion/deletion/substitution mutations by manual inspection and summarized the  
456 mutation as follows. The mutation length (d) was calculated as the difference between  
457 the length of the Sanger sequenced vector and the reference sequence, restricted to a  
458 window defined by high Sanger quality. The frameshift status was defined as (d) modulo  
459 3. To identify the indel location and length, we focused on the region between the  
460 translational start site and the mNeonGreen coding sequence. We then identified the  
461 first (reporter.prefix) and last (reporter.suffix) bases of the prefix and suffix sequences of  
462 the edited retrieval vector and the first (reference.prefix) and last (reference.suffix)  
463 bases of the prefix and suffix sequences of the corresponding region of the reference  
464 locus. We then defined 'query gap' and 'reference gap' as the difference between the  
465 prefix and suffix bases of the edited retrieval vector and the reference locus,  
466 respectively. (query gap = reporter.suffix - reporter.prefix; reference gap =  
467 reference.suffix - reference.prefix). The overall indel outcome was considered an  
468 insertion if the query gap exceeded the reference gap; otherwise, it was considered a  
469 deletion.

470

471 **Cell line authentication**

472 HeLa-TetR-Cas9 cells were a gift from Iain Cheeseman (MIT, Whitehead Institute).  
473 D458 cell-lines were a gift from Dr. Bigner (Duke University). To ensure the authenticity  
474 of cell lines, we performed Fluidigm SNP-based fingerprinting of each model cell line  
475 prior to screening. Cells were routinely tested to exclude the presence of mycoplasma.

476

477 **Declaration**

478 **Availability of data and materials**

479 The barcode read counts table for Figure 2 and Figure 4 are available in Additional file 4  
480 – Table S5-barcode\_counts.csv. Python scripts used for NGS analysis are available in  
481 Additional file 5 – Barcode count dataframe.ipynb and Additional file 7 – paella Python  
482 module. The raw histograms for barcode counts are available in Additional file 6 –  
483 barcode\_histograms. FCS files containing flow cytometry data supporting the  
484 conclusions of Figure 3 are available in Additional file 8 – Figure 3. Sanger sequencing  
485 datasets supporting the conclusions of Fig S4 are available in Additional file 9 – Sanger  
486 sequencing files.

487



488 **Competing interests**

489 R.B. and P.B. receive grant funding from the Novartis Institute of Biomedical Research  
490 for an unrelated project. C.M.J. is currently a full-time employee and stockholder of  
491 Novartis Institutes of BioMedical Research, Inc. The Broad Institute, Dana-Farber  
492 Cancer Institute, and MIT may seek to commercialize aspects of this work, and related  
493 applications for intellectual property have been filed.

494

495 **Funding**

496 This work was supported by SPARC funding from the Broad Institute (R.B. and C.M.J.),  
497 a grant from the Bridge Project of the Koch Institute for Integrative Cancer Research at  
498 MIT and the Dana-Farber/Harvard Cancer Center (R.B. and P.C.B.), the St Baldricks  
499 Foundation (P.B.), the Pediatric Brain Tumor Foundation (P.B. and R.B.), Alex's  
500 Lemonade Stand Foundation (R.B.), an NIH K99 award CA201592-02 (P.B.), NIH  
501 1U54CA224068-01 (C.M.J.), NIH R01 awards CA188228 (R.B.), CA219943 (R.B. and  
502 C.M.J.), and HG009283 (P.C.B.), the Jared Branfman Sunflowers for Life Fund for  
503 Pediatric Brain and Spinal Cancer Research (P.B. and R.B.). P.C.B. is supported by a  
504 Career Award at the Scientific Interface from the Burroughs Wellcome Fund.

505

## 506 **Author contributions**

507 D.F. and A.J.G. developed the frameshift retrieval vectors. D.F., F.T., A.J.G performed  
508 clone retrieval and characterization. D.F., F.T., A.J.G., R.O.R., L.B., P.H., E.G., and  
509 P.B. performed experiments. D.F. and F.T. analyzed data. P.C.B., C.M.J., R.B., and  
510 P.B. supervised the research. D.F., F.T., P.C.B., C.M.J., R.B., and P.B. wrote the  
511 manuscript with contributions from all authors.

512

## 513 **Acknowledgements**

514 We thank Emily Botelho and members of Blainey lab for feedback and discussions and  
515 acknowledge the Broad Flow Cytometry Facility for experimental assistance.

516

## 517 **Additional files**

518 Additional file 1: Fig S1-S5 and Table S1 and S2 (DOCX, 1.4 MB),

519 Additional file 2: Table S3 (XLSX, 8.8 MB),

520 Additional file 3: Table S4 (XLSX, 299.8 KB)

521 Additional file 4: Table S5-barcode\_counts (CSV, 7.6 MB),

522 Additional file 5: Barcode count dataframe (IPYNB, 5.3 KB).

523 Additional file 6: [barcode\\_histograms \(HIST.zip, 5.9 MB\)](#)

524 Additional file 7: [Python paella module \(PY, 131 KB\)](#)

525 Additional file 8: [Figure 3 \(FCS, 1.82 GB\)](#)

526 Additional file 9: [Sanger sequencing files \(SEQ, 393 KB\)](#)

527

## 528 **Figure Legends**

529 **Figure 1. Overview of the strategy for tracking and retrieving the ancestral clones**  
530 **within a heterogeneous population.**

531 (a) Tracking clonal response to selection (e.g.,  $\pm$ drug) using a lentiviral sgRNA-barcode  
532 library. Clonal fitness profiles can be estimated from barcode enrichment across  
533 replicates within each condition. (b) Clones of interest may be retrieved from the  
534 ancestral (untreated) population using a retrieval vector containing a targeting region  
535 matched to the clone sgRNA-barcode. Nuclease activity at the target region activates a  
536 fluorescent marker that can be detected with FACS. (c) Diagram of the frameshift  
537 retrieval vector. In cells from the clone of interest, where the sgRNA-barcode and  
538 barcode targets are matched, Cas9-mediated cleavage can induce a -1/+2 frameshift,  
539 activating reporter expression and inactivating mCherry expression. GFP+/mCherry-  
540 cells can be isolated by FACS. Additional reporter genes enable pre-enrichment such  
541 as antibiotic selection (e.g., zeocin) or affinity selection (e.g., H2K surface epitope) prior  
542 to FACS.

543

544 **Figure 2. Tracking clonal dynamics in D458 cells using a 26nt sgRNA-barcode**  
545 **library.**

546 (a) Relative barcode abundance in D458 cells before treatment (early time point, ETP)  
547 and after treatment with 2  $\mu$ M JQ1 (5 replicates) or DMSO vehicle (5 replicates). (b & c)  
548 The sgRNA-barcode library is able to track a heritable phenotype. (b) Comparison of  
549 barcode abundance across conditions for barcodes enriched in JQ1, DMSO, or JQ1  
550 and DMSO replicates. Barcode enrichment was defined based on the median rank  
551 across replicates (Methods). (c) The majority of JQ1-enriched barcodes were detected  
552 across all replicates at an abundance  $>10^{-5}$ . The raw barcode read counts are provided  
553 as a CSV file in Additional file 4: Table S5-barcode\_counts.csv and the raw histograms  
554 for barcode counts are available in Additional file 6 – barcode\_histograms.

555

556 **Figure 3. Retrieval vector performance.**

557 (a) HeLa cells were transduced with individual sgRNA-barcodes and paired with  
558 matched or mismatched barcode targets. Cells with frameshift +2 and 0 are expected to  
559 express GFP and mCherry, respectively, whereas cells with a +1 frameshift should  
560 express neither. (b) FACS analysis of plots of the TMv2 and TMv2-Zeo retrieval vectors  
561 with matched or mismatched barcode targets. (c) Incorporating tandem targets into the  
562 retrieval vectors enables multiplexed activation of a single vector by several barcodes.  
563 Gating strategy for analysis of the frameshift status of the cells is shown in Additional  
564 file 1: Fig. S5. The source data are provided as FCS files in Additional file 8.

565 **Figure 4. Retrieval of hygromycin-resistant clones from a heterogeneous**  
566 **population of HeLa cells.**

567 (a) Workflow to identify resistant clones using an sgRNA-barcode library. (Barcoding) A  
568 mixed population of hygromycin-resistant and hygromycin-sensitive HeLa cells was  
569 transduced with sgRNA-barcodes. (Selection) The resulting library was bottlenecked to  
570 limit barcode complexity, re-expanded, and cryo-preserved to define an early time point  
571 (ETP). Cells were then treated with either hygromycin or vehicle control (PBS).  
572 Hygromycin-enriched barcodes were determined by NGS. (b) Hygromycin-resistant  
573 barcodes were enriched across hygromycin-treated replicates. Barcode abundance for  
574 T1 (hygromycin-sensitive barcode candidate), T2 (hygromycin-resistant barcode  
575 candidate) and T3 (hygromycin-resistant barcode candidate). The raw barcode read  
576 counts are provided as a CSV file in Additional file 4: Table S5-barcode\_counts.csv and  
577 the raw histograms for barcode counts are available in Additional file 6 –  
578 barcode\_histograms (c) Workflow to retrieve resistant clones using the frameshift  
579 reporter. (Retrieval vector transduction) Hygromycin-sensitive and resistant candidate  
580 barcodes were selected for retrieval, and the matching barcode targets were cloned into  
581 the retrieval vector. Cells from the ETP were transduced with barcode-specific retrieval  
582 vectors and Cas9 expression was induced. (Clone enrichment and isolation) FACS  
583 sorting or zeocin selection was used to enrich for barcodes of interest. Single-cell  
584 clones were isolated by FACS. (Characterization) Barcode identification and functional  
585 validation. The integrated retrieval vector was sequenced to characterize specific and  
586 nonspecific mutations leading to reporter activation. (d) The ETP abundance of each

587 targeted barcode. (e) Population-level enrichment of targeted barcodes using selection  
588 by FACS (TMv2) or Zeocin selection (TMv2-Zeo). (f) Fraction of single-cell clones with  
589 the targeted barcode. (g) The hygromycin sensitivity of single-cell clones isolated by  
590 FACS corresponded to the sensitivity predicted by clonal tracking.

591

## 592 **References**

- 593 [1] S. Bonhoeffer and M. A. Nowak, "Pre-existence and emergence of drug  
594 resistance in HIV-1 infection.," *Proc. Biol. Sci.*, vol. 264, no. 1382, pp. 631–637,  
595 May 1997.
- 596 [2] F. Michor, T. P. Hughes, Y. Iwasa, S. Branford, N. P. Shah, C. L. Sawyers, and  
597 M. A. Nowak, "Dynamics of chronic myeloid leukaemia.," *Nature*, vol. 435, no.  
598 7046, pp. 1267–1270, Jun. 2005.
- 599 [3] M. Gerlinger, A. J. Rowan, S. Horswell, M. Math, J. Larkin, D. Endesfelder, E.  
600 Gronroos, P. Martinez, N. Matthews, A. Stewart, P. Tarpey, I. Varela, B.  
601 Phillimore, S. Begum, N. Q. McDonald, A. Butler, D. Jones, K. Raine, C.  
602 Latimer, C. R. Santos, M. Nohadani, A. C. Eklund, B. Spencer-Dene, G. Clark,  
603 L. Pickering, G. Stamp, M. Gore, Z. Szallasi, J. Downward, P. A. Futreal, and C.  
604 Swanton, "Intratumor heterogeneity and branched evolution revealed by  
605 multiregion sequencing.," *N. Engl. J. Med.*, vol. 366, no. 10, pp. 883–892, Mar.  
606 2012.
- 607 [4] W. J. Gibson, E. A. Hoivik, M. K. Halle, A. Taylor-Weiner, A. D. Cherniack, A.  
608 Berg, F. Holst, T. I. Zack, H. M. J. Werner, K. M. Staby, M. Rosenberg, I. M.  
609 Stefansson, K. Kusonmano, A. Chevalier, K. K. Mauland, J. Trovik, C. Krakstad,  
610 M. Giannakis, E. Hodis, K. Woie, L. Bjorge, O. K. Vintermyr, J. A. Wala, M. S.  
611 Lawrence, G. Getz, S. L. Carter, R. Beroukhim, and H. B. Salvesen, "The  
612 genomic landscape and evolution of endometrial carcinoma progression and  
613 abdominopelvic metastasis.," *Nat. Genet.*, vol. 48, no. 8, pp. 848–855, Aug.  
614 2016.
- 615 [5] P. Bandopadhyay, F. Piccioni, R. O'Rourke, P. Ho, E. M. Gonzalez, G. Buchan,  
616 K. Qian, G. Gionet, E. Girard, M. Coxon, M. G. Rees, L. Brenan, F. Dubois, O.  
617 Shapira, N. F. Greenwald, M. Pages, A. Balboni Iniguez, B. R. Paoletta, A.  
618 Meng, C. Sinai, G. Roti, N. V. Dharia, A. Creech, B. Tanenbaum, P. Khadka, A.  
619 Tracy, H. L. Tiv, A. L. Hong, S. Coy, R. Rashid, J.-R. Lin, G. S. Cowley, F. C.  
620 Lam, A. Goodale, Y. Lee, K. Schoolcraft, F. Vazquez, W. C. Hahn, A. Tsherniak,  
621 J. E. Bradner, M. B. Yaffe, T. Milde, S. M. Pfister, J. Qi, M. Schenone, S. A.  
622 Carr, K. L. Ligon, M. W. Kieran, S. Santagata, J. M. Olson, P. C. Gokhale, J. D.

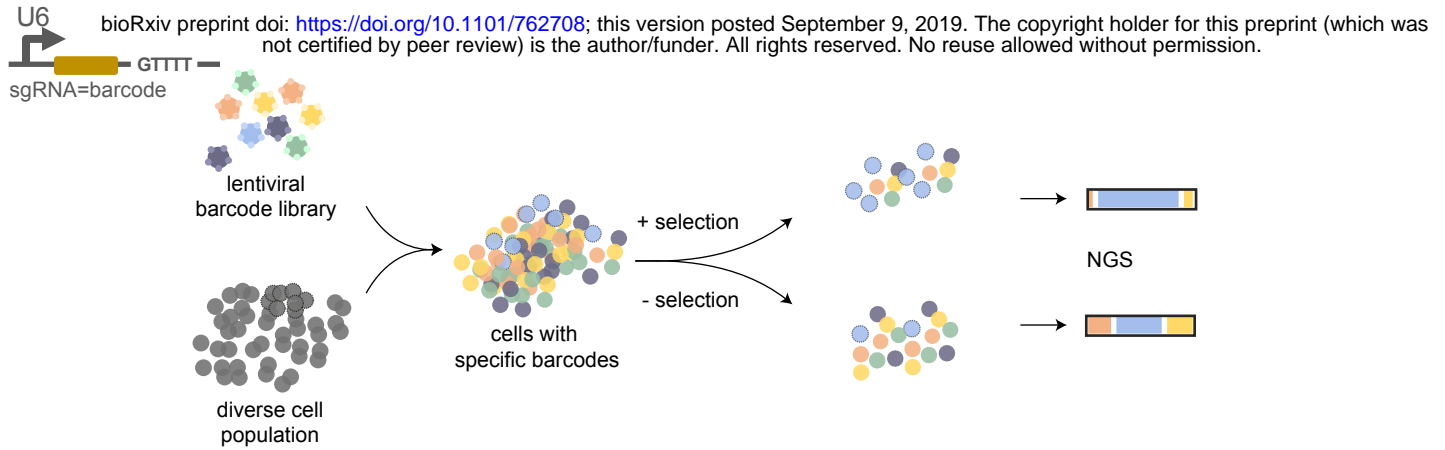
- 623 Jaffe, D. E. Root, K. Stegmaier, C. M. Johannessen, and R. Beroukhim,  
624 “Neuronal differentiation and cell-cycle programs mediate response to BET-  
625 bromodomain inhibition in MYC-driven medulloblastoma.,” *Nat Commun*, vol. 10,  
626 no. 1, p. 2400, Jun. 2019.
- 627 [6] J. Yang, S. A. Mani, J. L. Donaher, S. Ramaswamy, R. A. Itzykson, C. Come, P.  
628 Savagner, I. Gitelman, A. Richardson, and R. A. Weinberg, “Twist, a master  
629 regulator of morphogenesis, plays an essential role in tumor metastasis.,” *Cell*,  
630 vol. 117, no. 7, pp. 927–939, Jun. 2004.
- 631 [7] D. X. Nguyen, P. D. Bos, and J. Massagué, “Metastasis: from dissemination to  
632 organ-specific colonization.,” *Nat. Rev. Cancer*, vol. 9, no. 4, pp. 274–284, Apr.  
633 2009.
- 634 [8] H.-E. C. Bhang, D. A. Ruddy, V. Krishnamurthy Radhakrishna, J. X. Caushi, R.  
635 Zhao, M. M. Hims, A. P. Singh, I. Kao, D. Rakiec, P. Shaw, M. Balak, A. Raza,  
636 E. Ackley, N. Keen, M. R. Schlabach, M. Palmer, R. J. Leary, D. Y. Chiang, W.  
637 R. Sellers, F. Michor, V. G. Cooke, J. M. Korn, and F. Stegmeier, “Studying  
638 clonal dynamics in response to cancer therapy using high-complexity  
639 barcoding.,” *Nat. Med.*, vol. 21, no. 5, pp. 440–448, May 2015.
- 640 [9] B. A. Bidy, W. Kong, K. Kamimoto, C. Guo, S. E. Waye, T. Sun, and S. A.  
641 Morris, “Single-cell mapping of lineage and identity in direct reprogramming.,”  
642 *Nature*, vol. 564, no. 7735, pp. 219–224, Dec. 2018.
- 643 [10] S. Konermann, D. Feldman, F. Zhang, P. BLAINEY, The Broad Institute Inc.,  
644 Massachusetts Institute Of Technology, “Cell sorting,” *Plant Methods*, vol. 9, no.  
645 1, p. 39, Dec. 2016.
- 646 [11] M. T. Certo, B. Y. Ryu, J. E. Annis, M. Garibov, J. Jarjour, D. J. Rawlings, and A.  
647 M. Scharenberg, “Tracking genome engineering outcome at individual DNA  
648 breakpoints.,” *Nat. Methods*, vol. 8, no. 8, pp. 671–676, Jul. 2011.
- 649 [12] S. D. Perli, C. H. Cui, and T. K. Lu, “Continuous genetic recording with self-  
650 targeting CRISPR-Cas in human cells.,” *Science*, vol. 353, no. 6304, pp.  
651 aag0511–aag0511, Sep. 2016.
- 652 [13] U. Ben-David, B. Siranosian, G. Ha, H. Tang, Y. Oren, K. Hinohara, C. A.  
653 Strathdee, J. Dempster, N. J. Lyons, R. Burns, A. Nag, G. Kugener, B. Cimini, P.  
654 Tsvetkov, Y. E. Maruvka, R. O'Rourke, A. Garrity, A. A. Tubelli, P.  
655 Bandopadhyay, A. Tsherniak, F. Vazquez, B. Wong, C. Birger, M. Ghandi, A.  
656 R. Thorner, J. A. Bittker, M. Meyerson, G. Getz, R. Beroukhim, and T. R. Golub,  
657 “Genetic and transcriptional evolution alters cancer cell line drug response.,”  
658 *Nature*, vol. 560, no. 7718, pp. 325–330, Aug. 2018.
- 659 [14] H. Kim, E. Um, S.-R. Cho, C. Jung, H. Kim, and J.-S. Kim, “Surrogate reporters  
660 for enrichment of cells with nuclease-induced mutations.,” *Nat. Methods*, vol. 8,  
661 no. 11, pp. 941–943, Oct. 2011.
- 662 [15] S. Konermann, M. D. Brigham, A. E. Trevino, J. Joung, O. O. Abudayyeh, C.  
663 Barcena, P. D. Hsu, N. Habib, J. S. Gootenberg, H. Nishimasu, O. Nureki, and  
664 F. Zhang, “Genome-scale transcriptional activation by an engineered CRISPR-  
665 Cas9 complex,” *Nature*, vol. 517, no. 7536, pp. 583–588, Jan. 2015.
- 666 [16] M. van Overbeek, D. Capurso, M. M. Carter, M. S. Thompson, E. Frias, C. Russ,

- 667 J. S. Reece-Hoyes, C. Nye, S. Gradia, B. Vidal, J. Zheng, G. R. Hoffman, C. K.  
668 Fuller, and A. P. May, “DNA Repair Profiling Reveals Nonrandom Outcomes at  
669 Cas9-Mediated Breaks.,” *Mol. Cell*, vol. 63, no. 4, pp. 633–646, Aug. 2016.
- 670 [17] W. An and A. Telesnitsky, “Frequency of direct repeat deletion in a human  
671 immunodeficiency virus type 1 vector during reverse transcription in human  
672 cells.,” *Virology*, vol. 286, no. 2, pp. 475–482, Aug. 2001.
- 673 [18] A. Dixit, O. Parnas, B. Li, J. Chen, C. P. Fulco, L. Jerby-Arnon, N. D. Marjanovic,  
674 D. Dionne, T. Burks, R. Raychowdhury, B. Adamson, T. M. Norman, E. S.  
675 Lander, J. S. Weissman, N. Friedman, and A. Regev, “Perturb-Seq: Dissecting  
676 Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic  
677 Screens.,” *Cell*, vol. 167, no. 7, pp. 1853–1866.e17, Dec. 2016.
- 678 [19] P. Datlinger, A. F. Rendeiro, C. Schmidl, T. Krausgruber, P. Traxler, J.  
679 Klughammer, L. C. Schuster, A. Kuchler, D. Alpar, and C. Bock, “Pooled  
680 CRISPR screening with single-cell transcriptome readout.,” *Nat. Methods*, vol.  
681 14, no. 3, pp. 297–301, Mar. 2017.
- 682 [20] D. Feldman, A. Singh, J. L. Schmid-Burgk, A. Mezger, A. J. Garrity, R. J.  
683 Carlson, F. Zhang, and P. C. Blainey, “Pooled optical screens in human cells,”  
684 *bioRxiv*, p. 383943, Aug. 2018.
- 685 [21] “Cell Biology by the Numbers,” 2015.
- 686 [22] G. R. Abecasis, S. S. Cherny, W. O. Cookson, and L. R. Cardon, “Merlin—rapid  
687 analysis of dense genetic maps using sparse gene flow trees,” *Nat. Genet.*, vol.  
688 30, no. 1, pp. 97–101, Jan. 2002.
- 689 [23] T. 1. G. P. Consortium, “A map of human genome variation from population-  
690 scale sequencing,” *Nature*, vol. 467, no. 7319, pp. 1061–1073, Oct. 2010.
- 691 [24] M. W. Shen, M. Arbab, J. Y. Hsu, D. Worstell, S. J. Culbertson, O. Krabbe, C. A.  
692 Cassa, D. R. Liu, D. K. Gifford, and R. I. Sherwood, “Predictable and precise  
693 template-free CRISPR editing of pathogenic variants.,” *Nature*, vol. 563, no.  
694 7733, pp. 646–651, Nov. 2018.
- 695 [25] N. Shakiba, A. Fahmy, G. Jayakumaran, S. McGibbon, L. David, D. Trcka, J.  
696 Elbaz, M. C. Puri, A. Nagy, D. van der Kooy, S. Goyal, J. L. Wrana, and P. W.  
697 Zandstra, “Cell competition during reprogramming gives rise to dominant  
698 clones.,” *Science*, vol. 364, no. 6438, p. eaan0925, Apr. 2019.
- 699 [26] A. N. Hata, M. J. Niederst, H. L. Archibald, M. Gomez-Caraballo, F. M. Siddiqui,  
700 H. E. Mulvey, Y. E. Maruvka, F. Ji, H.-E. C. Bhang, V. Krishnamurthy  
701 Radhakrishna, G. Siravegna, H. Hu, S. Raoof, E. Lockerman, A. Kalsy, D. Lee,  
702 C. L. Keating, D. A. Ruddy, L. J. Damon, A. S. Crystal, C. Costa, Z. Piotrowska,  
703 A. Bardelli, A. J. Iafrate, R. I. Sadreyev, F. Stegmeier, G. Getz, L. V. Sequist, A.  
704 C. Faber, and J. A. Engelman, “Tumor cells can follow distinct evolutionary  
705 paths to become resistant to epidermal growth factor receptor inhibition.,” *Nat.*  
706 *Med.*, vol. 22, no. 3, pp. 262–269, Mar. 2016.
- 707 [27] N. E. Sanjana, O. Shalem, and F. Zhang, “Improved vectors and genome-wide  
708 libraries for CRISPR screening.,” *Nat. Methods*, vol. 11, no. 8, pp. 783–784,  
709 Aug. 2014.
- 710 [28] J. Joung, S. Konermann, J. S. Gootenberg, O. O. Abudayyeh, R. J. Platt, M. D.

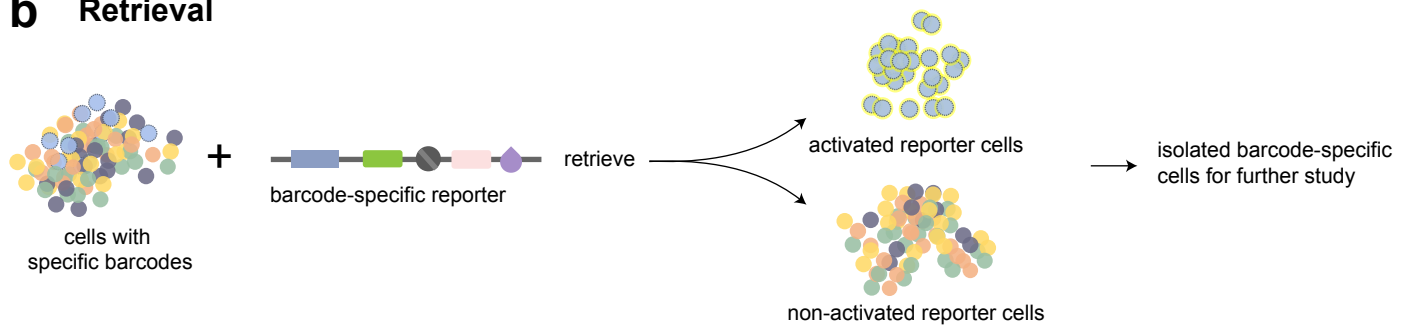


711 Brigham, N. E. Sanjana, and F. Zhang, “Genome-scale CRISPR-Cas9 knockout  
712 and transcriptional activation screening.,” *Nat Protoc*, vol. 12, no. 4, pp. 828–  
713 863, Apr. 2017.  
714

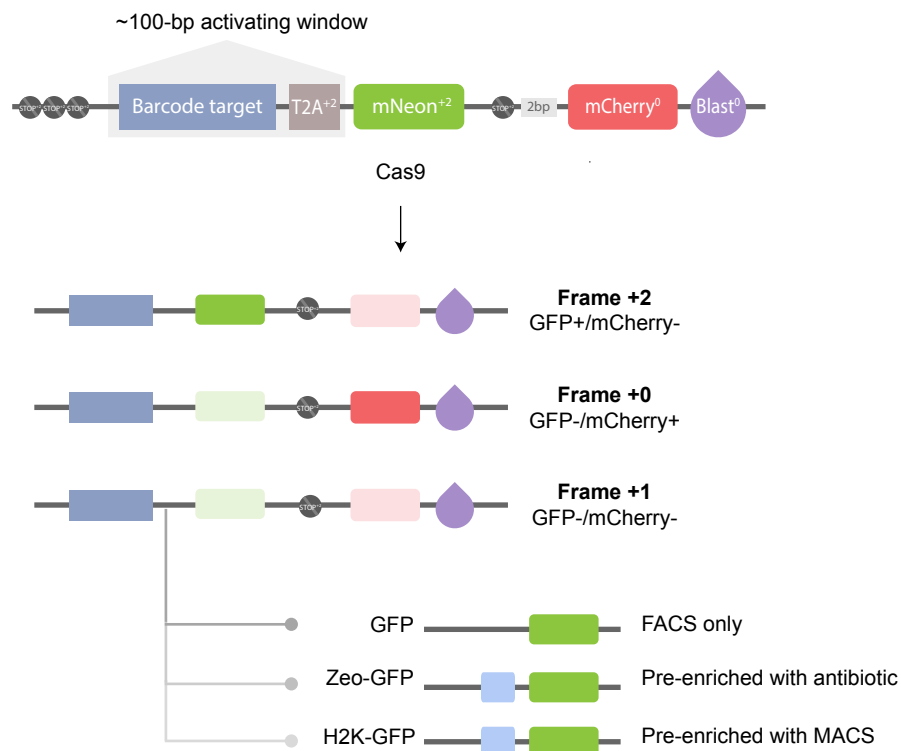
## a Clonal Tracking

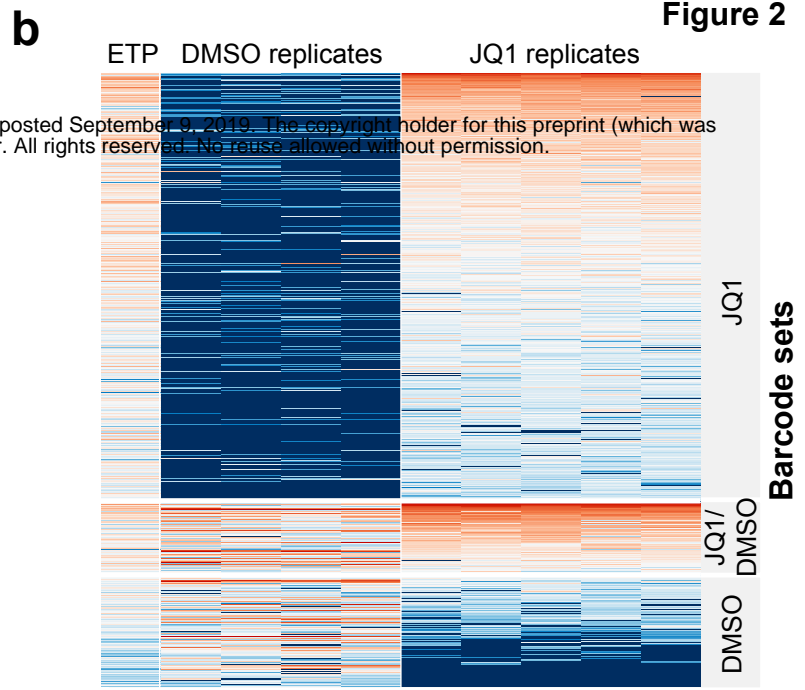
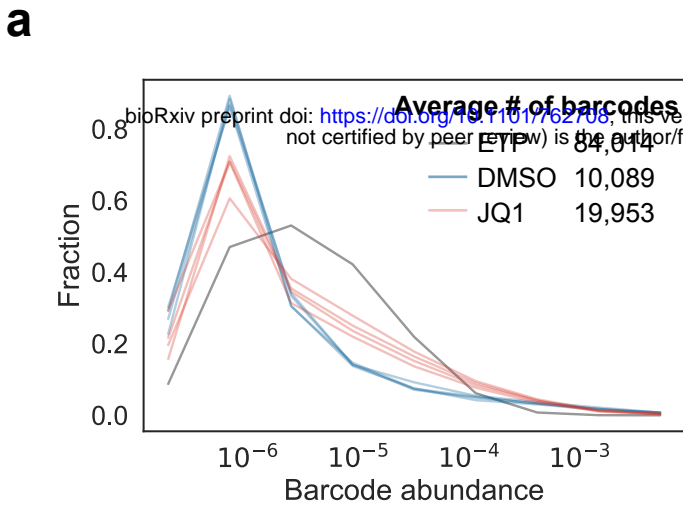


## b Retrieval



## c Retrieval Vector



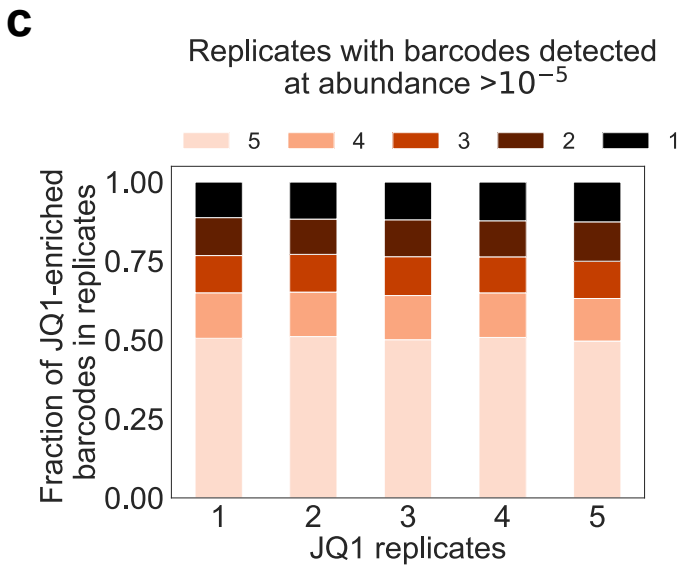


**# of enriched barcodes**

JQ1	2938	(2.39%)
JQ1 & DMSO	637	(0.52%)
DMSO	702	(0.57%)

Abundance (log10)

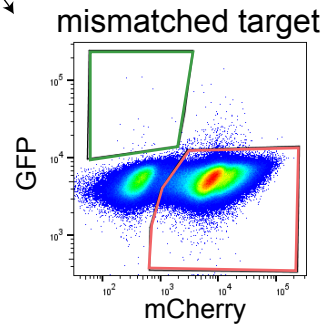
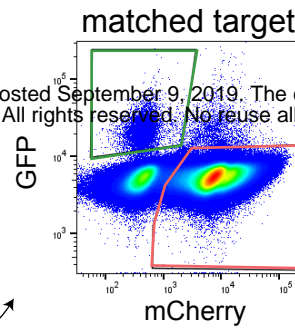
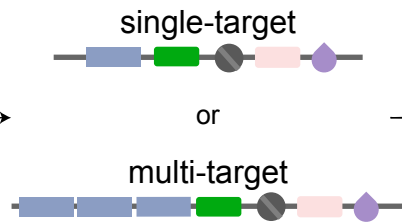
-1.5  
-3.0  
-4.5  
-6.0  
-7.5



**a**

bioRxiv preprint doi: <https://doi.org/10.1101/762708>; this version posted September 9, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

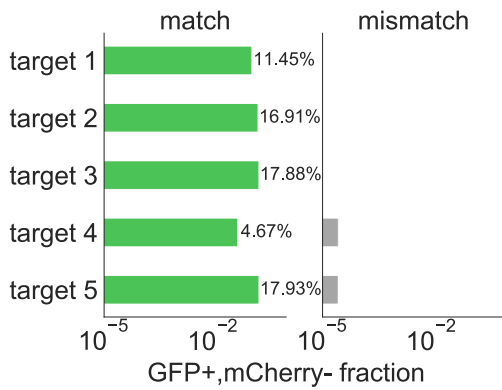
HeLa cells expressing sgRNA



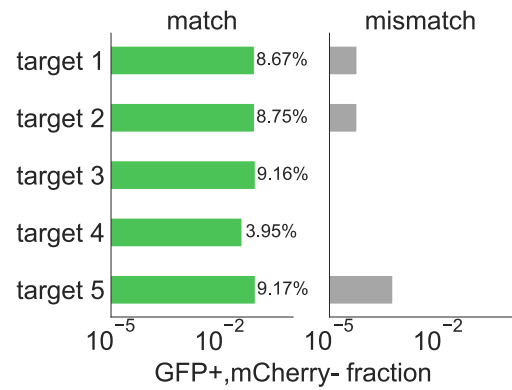
Reporter status	FACS outcome
+2	GFP+,mCherry-
+0	GFP-,mCherry+
+1	GFP-,mCherry-
low expression	GFP-,mCherry-

**b**

TMv2

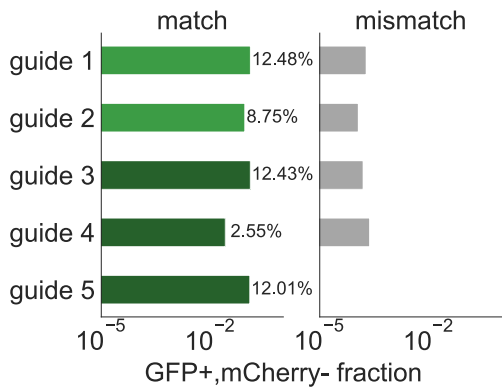


TMv2-Zeo

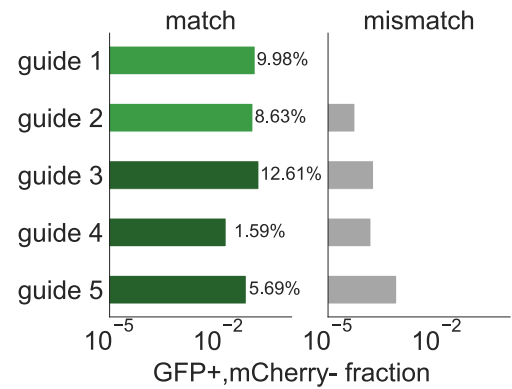


**c**

TMv2



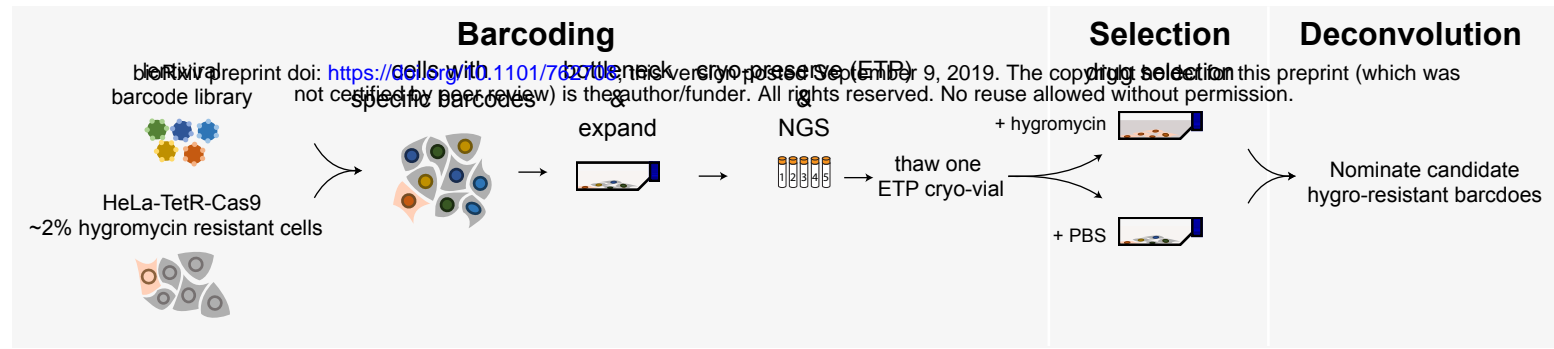
TMv2-Zeo



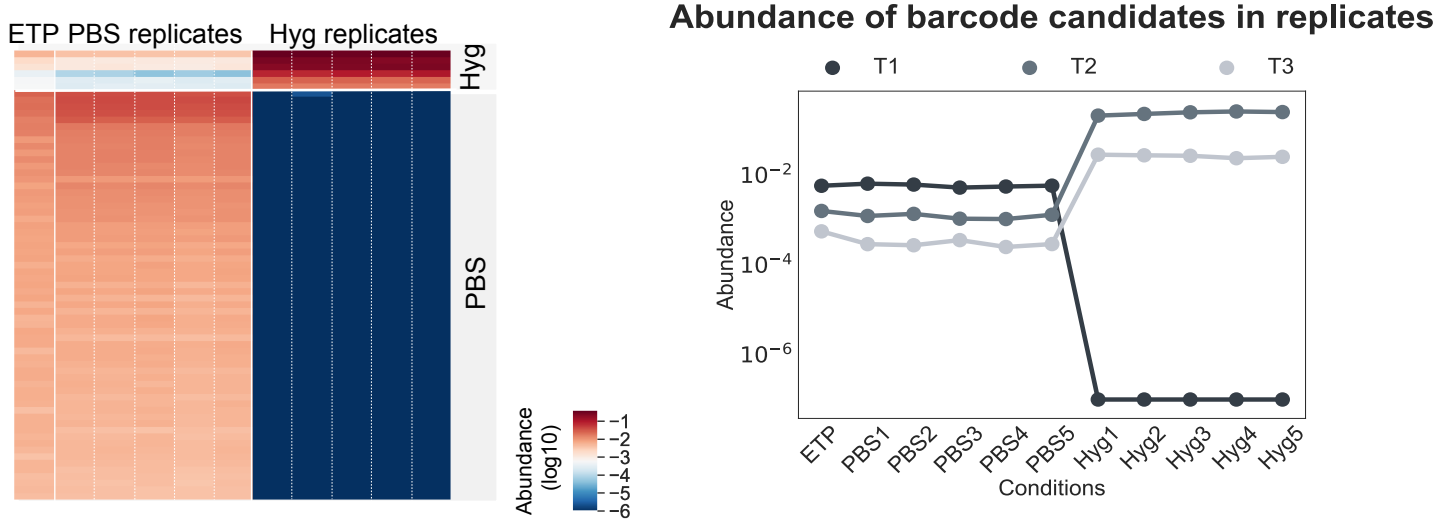
■ Multi-target 1

■ Multi-target 2

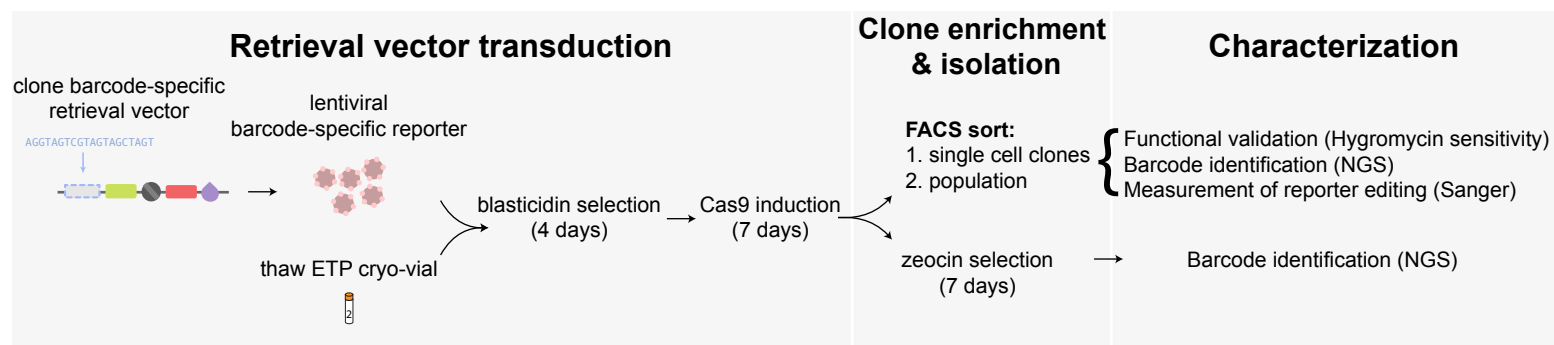
a



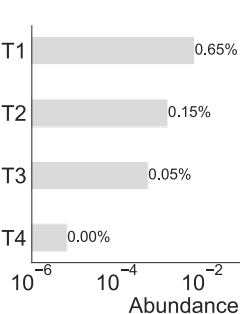
b



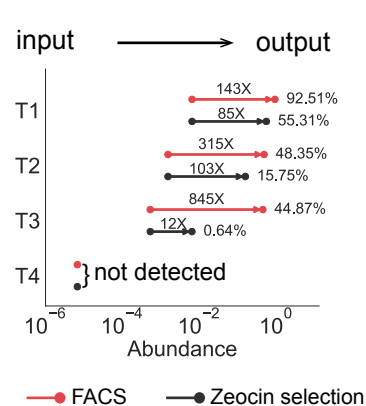
c



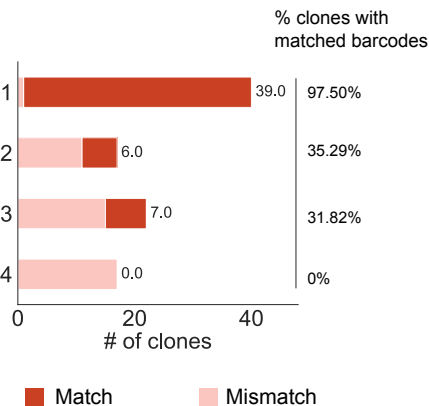
d



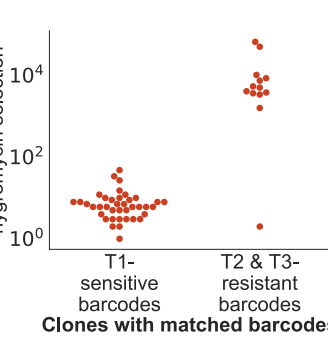
e



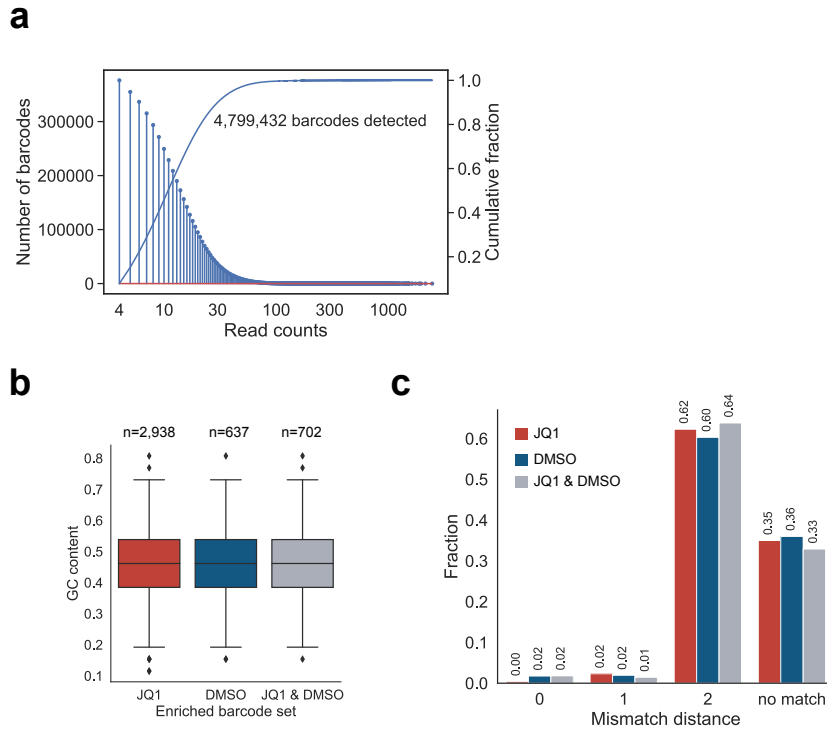
f



g



## 1 Supplementary Figures

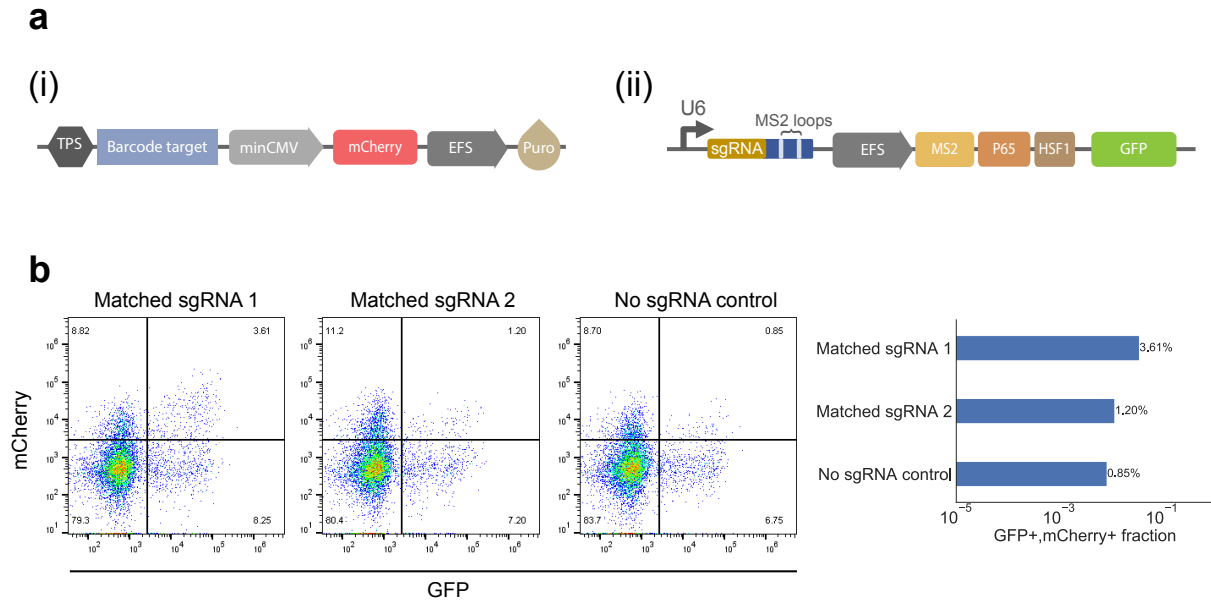


2 **Fig. S1 sgRNA-barcode library representation, GC bias and human genome off-**  
3 **targets.** (a) Deep sequencing of the 26-nt sgRNA-barcode library plasmid pool. (b) GC  
4 content in the 26-nt sgRNA-barcode sequence. There is no significant difference  
5 between the barcode sets. ( $p = 0.07$ , one-way ANOVA) (c) Distance of sgRNA-  
6 barcodes to human genome predicted by an off-target sgRNA algorithm [1]. The vast  
7 majority of sgRNA-barcodes have mismatch distance  $\geq 2$  homology to human genome.

8

9

10

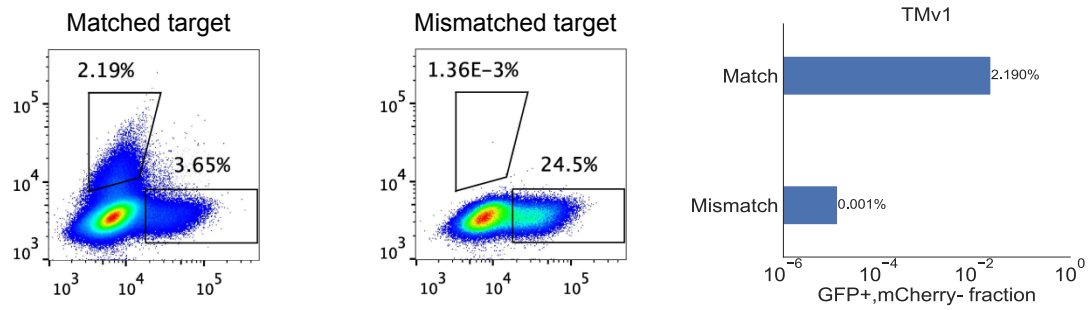


11

12 **Fig. S2 Transcriptional activation-based retrieval reporter.** (a) (i) The reporter  
13 comprises transcriptional pause site (TPS), barcode target, minCMV promoter, mCherry  
14 fluorescent reporter and puro resistance marker. (ii) sgRNA with transcriptional activator.  
15 (b) The reporter selectively induces expression of mCherry in cells with matching sgRNA  
16 (Matched sgRNA 1 & Matched sgRNA 2), while cells without the sgRNA sequence exhibit  
17 a low level of mCherry expression (No sgRNA control).

18

19



20

21 **Fig. S3 Specificity and sensitivity of the initial retrieval vector design (TMv1).**

22 TMv1 demonstrates high specificity, the background activation is around 1-4 in 100,000

23 cells (Matched sequence: GAGACCAGCAGAACCGACAA; Mismatched sequence:

24 GCGCAACAGAGAGGGGAGCG).

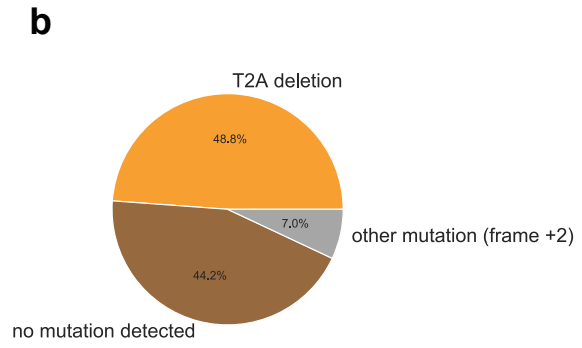
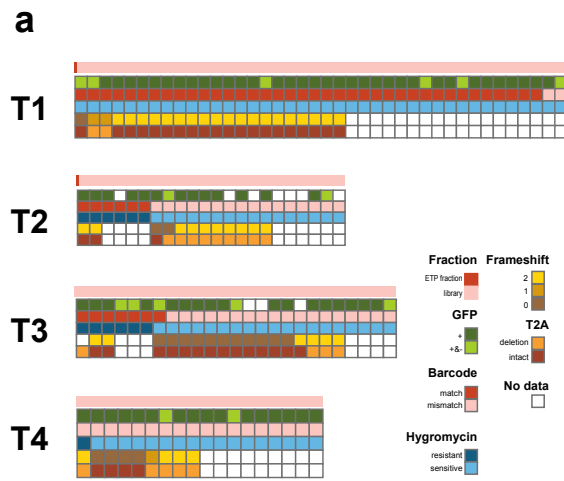
25

26

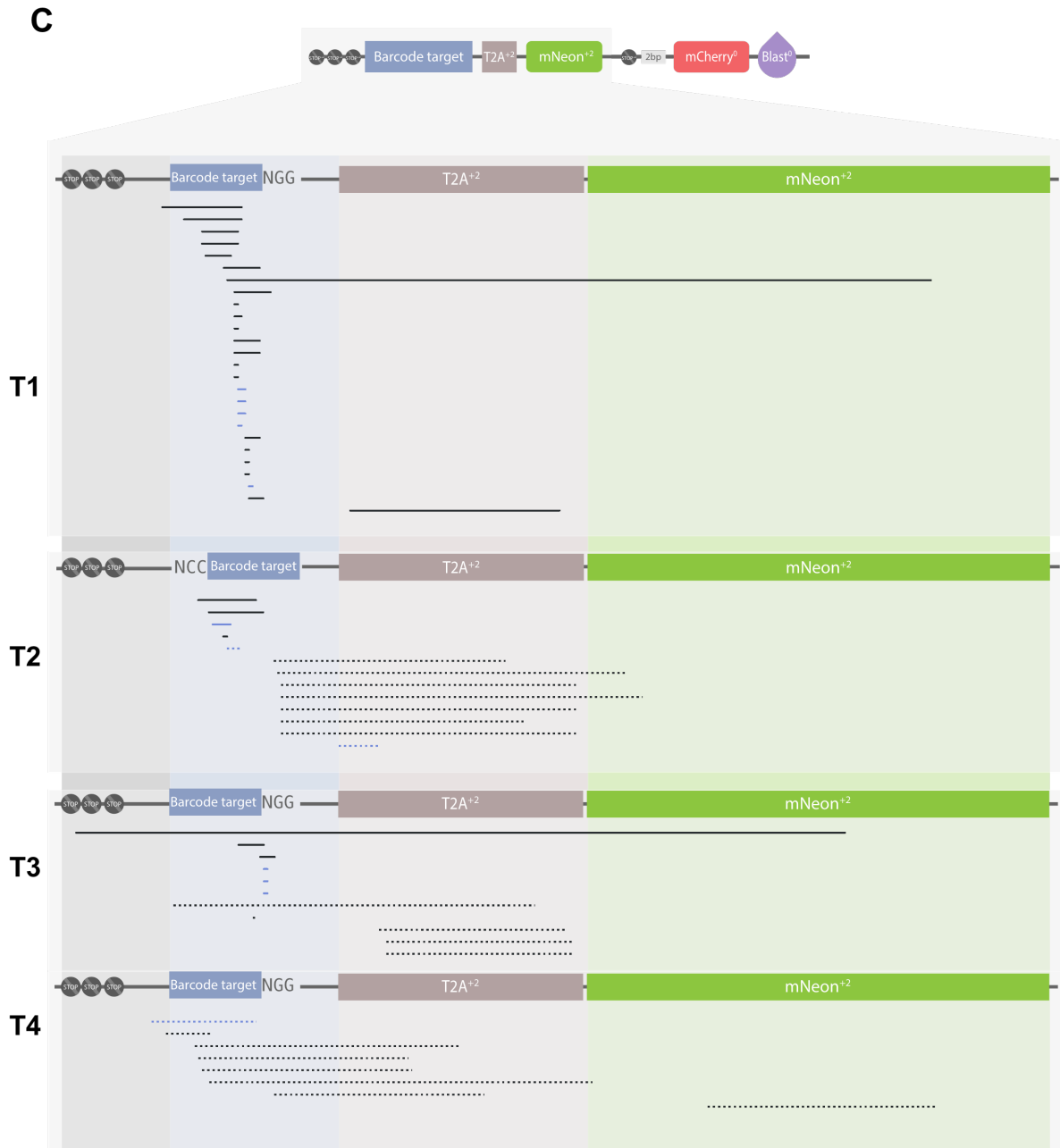
27

28





29



| clones with matched barcodes have deletion  
: clones with mismatched barcodes have deletion

| clones with matched barcodes have insertion  
: clones with mismatched barcodes have insertion

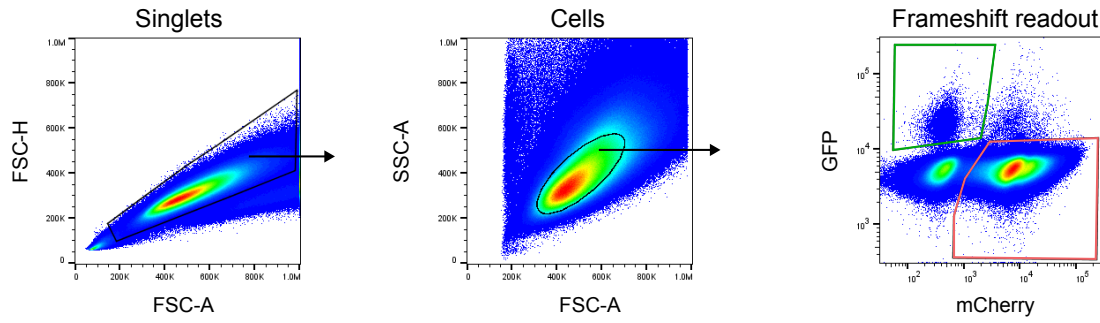
31 **Fig. S4 Characterization of retrieved clones.** (a and b) Instances of retrieval of clones  
32 carrying incorrect barcodes are primarily explained not by FACS error, but by  
33 background mutations, in particular a stereotyped T2A linker deletion. (a) Sorted clones  
34 were analyzed by FACS for GFP expression (green), validated for barcode accuracy by  
35 NGS (red), validated for hygromycin sensitivity (blue) and Sanger sequenced to  
36 determine frameshift status (brown) and T2A (orange) deletion. (b) The types of  
37 background in the clones with mismatched barcodes are shown. T2A deletions account  
38 for 48.8% of background events. (c) Map of the retrieval vector, focusing on the  
39 targeting region. Each line represents a clone, with black and blue lines represent  
40 deletion and insertion regions, respectively and normal and dashed lines representing  
41 clones with matched and mismatched barcodes, respectively. The line depicts the  
42 location of deletion and the length is proportional to the size of deletion and insertion.  
43 Sanger sequencing data for each clone is provided as a SEQ file in Additional file 9:  
44 Sanger sequencing files.

45

46

47

48



49

50 **Fig S5. Gating strategy for analysis of cells with activated frameshift reporter.**

51 Representative flow plots for HeLa cells after frameshift mutation induced with Cas9.

52

53

54

55

56

57

58

59

60

61

62 **Table S1. Table of clonal barcode candidates with their sgRNA-barcode sequences,**  
 63 **corresponding target sequences and hygromycin sensitivity.**

Clonal barcode	sgRNA-barcode	Target sequence (PAM lowercase)	Hygromycin sensitivity
T1	AGATCGTACCAGGGATT GGG	GCCACCATGCGCCAGATCG TACCAGGGATTGGG <b>agg</b> CTA CTGTACCGATCACT	sensitive
T2	TGCAGTGGCCGTTGAC AAAT	GCCACCATGGTACCGCGGG <b>cct</b> ATTTGTCAACGGCCACTG CACTACTATCACT	resistant
T3	GCGGGATCATTGCAATT ATA	GCCACCATGGTACCGCGCC GCGGGATCATTGCAATTATA <b>cgg</b> CTACTATCACT	resistant
T4	CAACATCCTGGGGCAC AAGC	GCCACCATGCGCCCAACAT CCTGGGGCACAAAGC <b>agg</b> CTA CTGTACCGATCACT	resistant

64  
65

66 **Table S2. Table of primer sequences used for amplifying 2 kb-lentiviral transgene**  
 67 **and for Sanger sequencing.**

Primers	Sequence (5'- 3')
pTM_negative_fwd	TCTTTCCCTACACGACGCTCTTCCGATCTAGCAGAG ATCCAGTTTGGTTAATTAGCTAGC
pTM_negative_rev	AAGACTACAGCGTCGCCAGCAGATCGGAAGAGCAC ACGTCTGAACTCCA
pTM_sanger_primer	GGATCTTGGTTCATTCTCAAGCC

68

69

70 **Table S1. Table of sgRNA-barcode sequences from the D458 clonal tracking**  
71 **experiment.** The abundance of each sgRNA-barcode was calculated from normalized  
72 read counts and transformed by a base-10 logarithm (Method). Median rank and median  
73 values of each sgRNA-barcode in each condition across replicates are listed. Identified  
74 barcode sets including JQ1, DMSO and JQ1 & DMSO are listed in the last column.

75

76 **Table S2. Table of sgRNA-barcode sequences from the HeLa clonal tracking**  
77 **experiment.** The abundance of each sgRNA-barcode was calculated with normalized  
78 read counts and transformed by a base-10 logarithm (Method). Median ranks and median  
79 values of each sgRNA-barcode in each condition across replicates are listed. Identified  
80 barcode sets including HeLa and PBS are listed in the last column.

81

## 82 **References**

83 [1] S. Bae, J. Park, and J.-S. Kim, “Cas-OFFinder: a fast and versatile algorithm that  
84 searches for potential off-target sites of Cas9 RNA-guided endonucleases.”  
85 *Bioinformatics*, vol. 30, no. 10, pp. 1473–1475, May 2014.