



# On the design of CRISPR-based single-cell molecular screens

Andrew J Hill<sup>1,3</sup> , José L McFaline-Figueroa<sup>1,3</sup>,  
Lea M Starita<sup>1</sup>, Molly J Gasperini<sup>1</sup>, Kenneth A Matreyek<sup>1</sup>,  
Jonathan Packer<sup>1</sup>, Dana Jackson<sup>1</sup>, Jay Shendure<sup>1,2,4</sup> &  
Cole Trapnell<sup>1,4</sup> 

**Several groups recently coupled CRISPR perturbations and single-cell RNA-seq for pooled genetic screens. We demonstrate that vector designs of these studies are susceptible to ~50% swapping of guide RNA–barcode associations because of lentiviral template switching. We optimized a published alternative, CROP-seq, in which the guide RNA also serves as the barcode, and here confirm that this strategy performs robustly and doubled the rate at which guides are assigned to cells to 94%.**

Pooled genetic screens based on RNAi or CRISPR enable thousands of programmed perturbations per experiment<sup>1,2</sup>. However, assays for such screens are limited to coarse phenotypes (e.g., cell viability) and are uninformative with respect to the mechanism by which perturbations mediate their effects.

To circumvent these limitations, several groups recently reported using single-cell RNA-seq (scRNA-seq) as a readout for CRISPR-based pooled genetic screens. The single guide RNA (sgRNA) in each cell is identified together with its transcriptome, either via a Pol II transcribed barcode (CRISP-seq, Perturb-seq, Mosaic-seq<sup>3–6</sup>) (Fig. 1a) or by capturing the sgRNA itself within a Pol II transcript (CROP-seq<sup>7</sup>) (Fig. 1b). Toward similar goals, we pursued a lentiviral strategy similar to former methods<sup>3–6</sup> in which each sgRNA was linked to a barcode located several kilobases away (Fig. 1a). In our vector (pLGB-scKO), the barcode was positioned in the 3' UTR of a blasticidin resistance transgene, enabling its recovery by scRNA-seq methods that capture poly(A) transcripts (Supplementary Fig. 1a,b). Guides and barcodes were paired during DNA synthesis, which facilitated pooled cloning and lentiviral delivery (Supplementary Fig. 1c).

With this design, we sought to ask how loss-of-function (LoF) of tumor suppressors altered gene expression in immortalized, non-transformed breast epithelial cells. We targeted *TP53* and other tumor suppressors in MCF10A cells, with or without exposure to the DNA-damaging agent doxorubicin. Cloning and lentiviral packaging was performed either individually for each targeted gene ('arrayed') or in a pooled fashion. In addition to scRNA-seq,

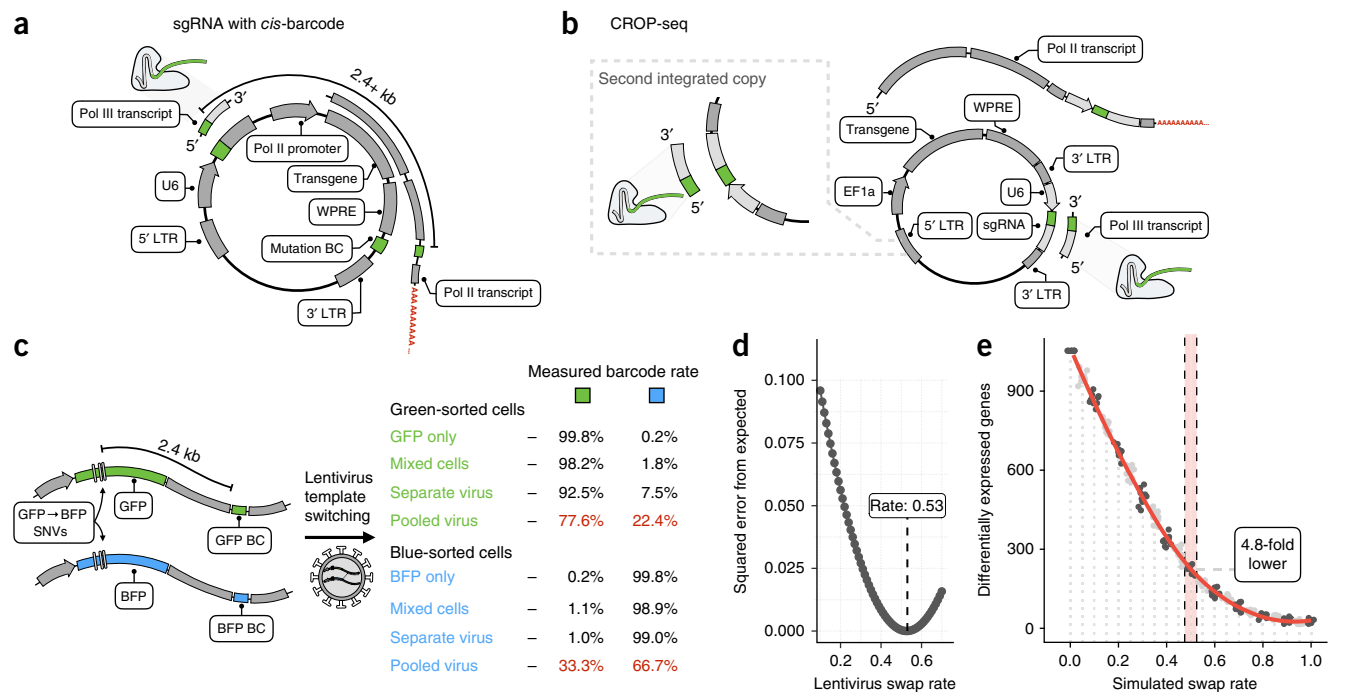
we performed targeted amplification<sup>4,5</sup> to more efficiently recover the barcodes present in each cell (Supplementary Fig. 1b and Supplementary Fig. 2).

With arrayed lentiviral production, a substantial proportion of cells in which *TP53* was targeted had a gene expression signature consistent with failure to activate a cell cycle checkpoint response after DNA damage (e.g., lower expression of *CDKN1A* and *TP53I3*; Supplementary Fig. 3a). However, these effects were greatly reduced when we performed a similar experiment with pooled lentiviral production (Supplementary Fig. 3b). Furthermore, markedly fewer genes were differentially expressed in the pooled than in the arrayed experiment (Supplementary Fig. 3c). t-SNE embedding revealed that both experiments contained a cluster of cells characterized by expression of the mitotic marker *CCNB2* and low levels of *TP53I3*, consistent with a *TP53*-null phenotype. In the arrayed experiment, this cluster was almost entirely composed of cells with sgRNAs targeting *TP53* (99.4%). However, in the pooled experiment, only 41% of assigned cells from the corresponding cluster contained *TP53* sgRNAs (Supplementary Fig. 3d–i).

We reasoned that lentiviral template switching may explain this difference. Lentiviral virions are pseudodiploid; i.e., two viral transcripts are copackaged during their production<sup>8,9</sup>. The reverse transcriptase that acts before integration has a rate of template switching<sup>10</sup> estimated as 1 event per kilobase (kb)<sup>11</sup>. In pooled lentiviral production, template switching should result in the integration of chimeric products at a rate proportional to the distance between paired sequences (Supplementary Fig. 4). This risk was noted by Adamson *et al.*<sup>4</sup> and Dixit *et al.*<sup>5</sup>. It was altogether avoided by Adamson *et al.*<sup>4</sup> through arrayed lentiviral production, but pooled lentiviral production was performed in some or all experiments of the other reports<sup>3,5,6</sup>. Although Sack *et al.*<sup>12</sup> recently quantified this phenomenon at distances up to 720 bp in vectors designed for bulk selection screens, the implications of template switching at longer distances (e.g., the 2.5 kb+ separation between sgRNAs and barcodes in the pLGB-scKO, CRISP-seq, Perturb-seq, and Mosaic-seq vectors<sup>3–6</sup>), as well as for scRNA-seq study designs specifically, remain unexplored.

To test this hypothesis, we cloned BFP and GFP transgenes, which differ by 3 bp, into separate lentiviral vectors, pairing each with a unique barcode separated from the nearest unique bases in BFP or GFP by 2.4 kb (Fig. 1c). We transduced MCF10A cells with lentivirus generated either individually or as a pool of the two plasmids, FACS-sorted GFP+ or BFP+ fractions, and we quantified the rate of barcode swapping (Fig. 1c and Supplementary Fig. 5). At this distance, swapping occurred at the theoretical maximum rate of 50% (Fig. 1d and Supplementary Fig. 6).

<sup>1</sup>Department of Genome Sciences, University of Washington, Seattle, Washington, USA. <sup>2</sup>Howard Hughes Medical Institute, Seattle, Washington, USA. <sup>3</sup>These authors contributed equally to this work. <sup>4</sup>These authors jointly directed this work. Correspondence should be addressed to C.T. (coletrap@uw.edu) or J.S. (shendure@uw.edu).

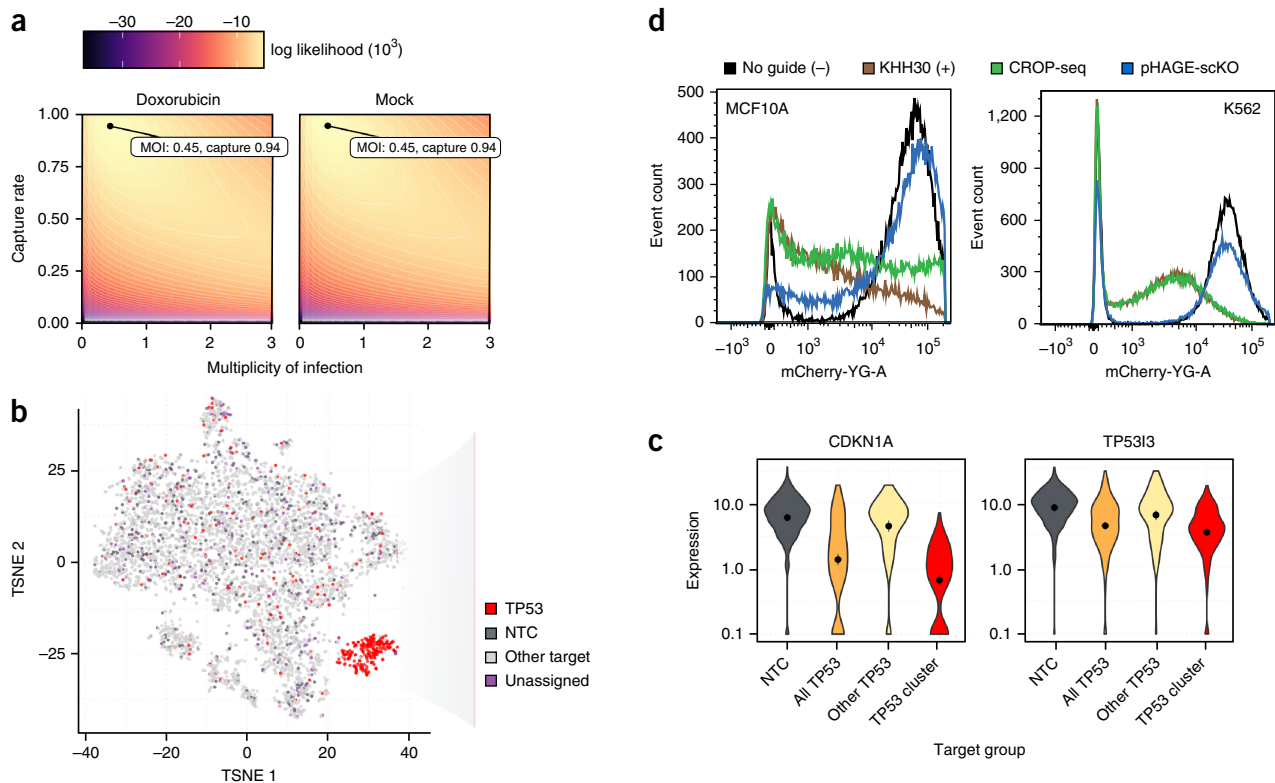


**Figure 1** | Template switching decreases the sensitivity of CRISPR-based single-cell molecular screens that employ linked barcodes. **(a)** Schematic of vectors that rely on *cis*-pairing of sgRNAs and barcodes such as Perturb-seq, CRISP-seq, and MOSAIC-seq. A barcode (BC), expressed as part of the Pol II transcript and sequenced as a proxy for the guide sequence, is linked to an sgRNA by a distance of 2.4 kb or more. WPRE, woodchuck hepatitis virus post-transcriptional regulatory element. U6, a Pol III promoter. **(b)** CROP-seq approach. One copy of the guide is cloned into the 3' LTR and transcribed as part of the Pol II transcript, which can be sequenced directly. A second copy of the guide expression cassette is produced in the 5' LTR during lentivirus positive-strand synthesis before integration. **(c)** Template switching at 2.4 kb separation between the distinguishing bases (3-bp differences) in GFP and BFP and their respective barcodes. Percentages reflect sorted cells transduced with GFP virus (GFP only) or BFP virus (BFP only); these cells mixed before sorting (mixed cells); or cells transduced with mixed virus generated from GFP and BFP plasmid packaged individually (separate virus) or together (pooled virus). Note that in a mix of two plasmids, only approximately half of all chimeric products are detectable because of homozygous virions (see Online Methods). **(d)** Sum of squared errors of observed data vs. expected values at various swap rates using the fraction of barcodes in the green and blue sorted samples ( $n = 4$  measurements), assuming a relative proportion of 61.7% GFP+ cells as determined from FACS (see **Supplementary Fig. 4** and Online Methods for details). **(e)** Simulation of progressively higher fractions of target assignment swapping on data from the transcription factor pilot arrayed screen of Adamson *et al.*<sup>4</sup>, used here as a gold standard performed with arrayed lentivirus production. Number of DEG across the target label at FDR of 5% is plotted at each swap rate for ten samplings per swap rate ( $n = 5,321$  cells used in tests). 0.5 corresponds to the 50% swap rate determined via FACS.

To simulate the impact of template switching, we obtained data from Adamson *et al.*<sup>4</sup> generated using the Perturb-seq vector with arrayed lentiviral production. We swapped target labels *in silico* at varying rates, and we evaluated power to detect differentially expressed genes (DEG). With 50% swapping, we observe a 4.8-fold decrease in the number of DEG (**Fig. 1e**). This loss in power results from an effective two-fold reduction in number of useful cells per target, coupled with noise from swapped associations.

CROP-seq<sup>7</sup> differs from the other methods<sup>3–6</sup> in that it does not rely on pairing of sgRNAs and barcodes. Instead, the sgRNA itself serves as a barcode as part of an overlapping Pol II transcript. Furthermore, the sgRNA cassette is copied from the 3' to 5' long terminal repeat (LTR) during positive-strand synthesis (**Fig. 1b**) via an intramolecular priming step that does not result in appreciable intermolecular swapping<sup>13</sup>. A limitation of CROP-seq is that sgRNAs are recovered from scRNA-seq data with limited sensitivity (~40–60%)<sup>7</sup>, such that half the single-cell transcriptomes are discarded. We modified CROP-seq to include targeted amplification of the sgRNA region from mRNA libraries already tagged with cellular barcodes, similar to our pLGB-scKO design (**Supplementary Fig. 7a,b**).

To evaluate this approach, we performed a CRISPR-mediated LoF screen of 32 tumor suppressors (six guides per target) and six nontargeting control (NTC) guides in MCF10A cells with or without doxorubicin. Whereas sgRNA(s) would generally be identified at a rate of 42–47% from scRNA-seq data alone, this rate was 94% with targeted amplification (**Fig. 2a**). In contrast with our original pooled experiment, tSNE embedding of doxorubicin-exposed cells from this experiment yielded a cluster almost entirely composed of cells containing *TP53*-targeting sgRNAs (**Fig. 2b**). Specifically, the 262 cells in this cluster include 90.5% with *TP53*-targeting guides, 7.6% with guides targeting other genes, 0% with NTC guides, and 1.9% unassigned cells. In contrast, the remaining 5,617 cells include 3.2% with *TP53*-targeting guides (presumably cells in which LoF editing failed to occur), 84.2% with guides targeting other genes, 7.5% with NTC guides, and 5.2% unassigned cells. Expression levels of the p53 targets *CDKN1A* and *TP53I3* (refs. 14 and 15) were markedly lower in the *TP53*-targeted cluster (**Fig. 2c**); and 4,277 and 2,186 DEGs (false discovery rate (FDR) 5%) were identified relative to cells with NTC guides in the doxorubicin-treated and untreated (mock) conditions, respectively. Thus, our improved CROP-seq protocol achieves the power and negligible sgRNA swap rate of the arrayed format without



**Figure 2** | CROP-seq with PCR enrichment offers improvements over alternate screen designs in a tumor suppressor knockout screen. **(a)** Most likely multiplicity of infection (MOI) and capture rate of barcoded transcripts in CROP-seq screen based on a generative model. **(b)** tSNE embedding of a doxorubicin-treated sample highlighting cells that contain *TP53* guides, nontargeting controls (NTC) or non-*TP53* guides, and unassigned cells ( $n = 5,879$  cells). **(c)** *CDKN1A* and *TP53I3* expression in cells expressing either nontargeting controls or *TP53* guides. Cells with *TP53* guides are further stratified by inclusion in the *TP53*-enriched cluster from **b**. Values below 0.1 are not shown to facilitate plotting. **(d)** CRISPRi knockdown of mCherry in MCF10A and K562 cells not expressing a guide (– control), KHH30 (+ control), CROP-seq, and pHAGE-scKO design. All vectors have been modified to contain a CRISPRi-optimized backbone. pHAGE-scKO places the sgRNA within a Pol II 3' UTR and does not knock down mCherry expression.

sacrificing the scalability of a pooled cloning and lentiviral production workflow.

Upon tSNE analysis of both mock and doxorubicin-treated cells (**Supplementary Fig. 8a,b**), we find several tumor suppressors whose distribution across clusters is significantly different compared to that of NTCs (FDR 5%; 13 and 14 targets with significant changes in the mock and doxorubicin conditions, respectively) (**Supplementary Fig. 8c–f**). We tested for target enrichment within clusters and generated average expression profiles for each enriched target–cluster pair. Gene set enrichment analysis of the most highly loaded genes in the principal components of these average expression profiles show many targets to be associated with increased proliferation and a decreased DNA damage response, most prominently with targeting of *TP53* (**Supplementary Fig. 9**).

To further assess the impact of template switching on sensitivity, we permuted target labels within our own CROP-seq tumor-suppressor screen, observing a 2.9-fold reduction in the number of DEGs across targets at a swap rate of 50%. The number of significant targets was also reduced, to just 4/13 (*TP53*, *STK11*, *CHEK1*, and *NCOR1*) and 3/14 (*TP53*, *RB1*, and *ARID1B*) in the mock and doxorubicin conditions, respectively. Additionally, simulations of 50% swapping on the larger (50,000 cells) unfolded-protein response screen from Adamson *et al.*<sup>4</sup> with arrayed lentiviral production resulted in a 1.9- and 2.8-fold reduction in the number of DEGs when using 25,000 and 6,000 cells, respectively

(**Supplementary Fig. 10**). Altogether, these simulations demonstrate that the reduction in power consequent to swapping is dependent on the number of cells captured, the number of targets, and the effect size of those targets.

Although CROP-seq is not subject to sgRNA-barcode swapping, it is limited by its placement of the sgRNA in the lentiviral LTR, as larger intervening sequences such as dual sgRNA designs<sup>16</sup> might render the LTR nonfunctional<sup>7</sup>. To enable incorporation of longer cassettes, we placed the sgRNA cassette between the WPRE and LTR. In this design (pHAGE-scKO), copying of the sgRNA between LTRs would not occur, but the guide sequence would still contribute to overlapping Pol II and Pol III transcripts (**Supplementary Fig. 11**).

To evaluate this design, we compared the ability of pHAGE-scKO, CROP-seq, and a standard lentiviral sgRNA expression vector, pKHH030 (ref. 17), all containing a CRISPRi-optimized backbone, to inhibit transcription via CRISPRi, targeting the promoter of an mCherry transgene. Whereas pKHH030 and CROP-seq exhibited efficient inhibition, pHAGE-scKO had poor efficacy (**Fig. 2d**). Consistent with this, we observed low editing rates with pHAGE-scKO (88% with pLGB control vs. 29% with pHAGE-scKO). Recent studies suggest interference when Pol II and Pol III transcripts overlap<sup>18,19</sup>. We hypothesize that the poor efficacy of pHAGE-scKO is due to the blasticidin resistance gene inhibiting sgRNA expression. In contrast, CROP-seq likely maintains

efficacy because the second integrated copy of the sgRNA (copied to the 5' LTR) does not overlap a Pol II transcript.

CRISPR-based pooled genetic screens coupled to scRNA-seq phenotyping have the potential to be extremely powerful. However, several published designs, and our own initial design, are susceptible to high rates of sgRNA-barcode swapping (diagrams of all relevant vectors are shown in **Supplementary Fig. 12**). Importantly, we do not expect that positive conclusions drawn by published studies using such designs in conjunction with pooled lentivirus production<sup>3,5,6</sup> are incorrect. Each of these studies examined few targets and collected large data sets, raising their baseline sensitivity. However, given the high cost of scRNA-seq and impetus to expand the number of targets in such screens, our observations are highly relevant for future studies. Reductions in power may be partly overcome by filtering cells that appear inconsistent with their assigned target<sup>5</sup>, or completely overcome with arrayed lentiviral production (as in Adamson *et al.*<sup>4</sup>). However, computational filtering has the potential to introduce biases, and itself reduces power by discarding collected data, while arrayed lentiviral production dramatically limits scalability.

A viable alternative is the recently published CROP-seq method<sup>7</sup>. By coupling targeted sgRNA amplification and CROP-seq, we doubled the proportion of cells in which guides are assigned to 94%. The attractive features of this approach include the simplicity of the cloning protocol, its compatibility with lentiviral delivery, the high rate of recovery of sgRNA-cell associations, and minimized risk of template switching.

## METHODS

Methods, including statements of data availability and any associated accession codes and references, are available in the [online version of the paper](#).

*Note: Any Supplementary Information and Source Data files are available in the online version of the paper.*

## ACKNOWLEDGMENTS

We thank all members of the Shendure and Trapnell labs for feedback on our manuscript and helpful discussions, particularly S. Srivatsan, G. Findlay, A. McKenna, R. Daza, B. Martin, M. Kircher, D. Cusanovich, X. Qiu, and V. Ramani. We thank J. Bloom and D. Fowler for discussions about lentivirus, and K. Han, J. Ousey, and M. Bassik for experimental advice and reagents for CRISPRi experiments. A.J.H. thanks Stella the cat for support. This work was supported by the following funding: NIH DP1HG007811 and UM1HG009408

(to J.S.), DP2HD088158 (to C.T.), and the W.M. Keck Foundation (to C.T. and J.S.). A.J.H. and M.J.G. are funded by the National Science Foundation Graduate Research Fellowship. J.L.M. is supported by the NIH Genome Training Grant (5T32HG000035) and the Cardiovascular Research Training Grant (4T32HL007828). C.T. is partly supported by an Alfred P. Sloan Foundation Research Fellowship. J.S. is an Investigator of the Howard Hughes Medical Institute.

## AUTHOR CONTRIBUTIONS

A.J.H., J.L.M.-F., J.S., and C.T. devised the project. A.J.H., J.L.M.-F., L.M.S., and M.J.G. performed experiments. D.J. optimized cloning strategies and provided substantial technical support. A.J.H., J.L.M.-F., and J.P. performed analysis. K.A.M. provided critical input on mechanisms of template switching in lentivirus. A.J.H., J.L.M., J.S. and C.T. wrote the manuscript with input from other authors.

## COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>. Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

- Shalem, O., Sanjana, N.E. & Zhang, F. *Nat. Rev. Genet.* **16**, 299–311 (2015).
- Mohr, S.E., Smith, J.A., Shamu, C.E., Neumüller, R.A. & Perrimon, N. *Nat. Rev. Mol. Cell Biol.* **15**, 591–600 (2014).
- Xie, S., Duan, J., Li, B., Zhou, P. & Hon, G.C. *Mol. Cell* **66**, 285–299.e5 (2017).
- Adamson, B. *et al. Cell* **167**, 1867–1882.e21 (2016).
- Dixit, A. *et al. Cell* **167**, 1853–1866.e17 (2016).
- Jaitin, D.A. *et al. Cell* **167**, 1883–1896.e15 (2016).
- Datlinger, P. *et al. Nat. Methods* **14**, 297–301 (2017).
- Nikolaitchik, O.A. *et al. PLoS Pathog.* **9**, e1003249 (2013).
- Tseng, W.C., Haselton, F.R. & Giorgio, T.D. *J. Biol. Chem.* **272**, 25641–25647 (1997).
- Jetzt, A.E. *et al. J. Virol.* **74**, 1234–1240 (2000).
- Schlub, T.E., Smyth, R.P., Grimm, A.J., Mak, J. & Davenport, M.P. *PLoS Comput. Biol.* **6**, e1000766 (2010).
- Sack, L.M., Davoli, T., Xu, Q., Li, M.Z. & Elledge, S.J. *G3 (Bethesda)* **6**, 2781–2790 (2016).
- Yu, H., Jetzt, A.E., Ron, Y., Preston, B.D. & Dougherty, J.P. *J. Biol. Chem.* **273**, 28384–28391 (1998).
- el-Deiry, W.S. *et al. Cell* **75**, 817–825 (1993).
- Contente, A., Dittmer, A., Koch, M.C., Roth, J. & Dobbstein, M. *Nat. Genet.* **30**, 315–320 (2002).
- Gasperini, M. *et al. Am. J. Hum. Genet.* **101**, 192–205 (2017).
- Han, K. *et al. Nat. Biotechnol.* **35**, 463–474 (2017).
- Lukoszek, R., Mueller-Roeber, B. & Ignatova, Z. *FEBS Lett.* **587**, 3692–3695 (2013).
- Yeganeh, M., Praz, V., Cousin, P. & Hernandez, N. *Genes Dev.* **31**, 413–421 (2017).

## ONLINE METHODS

**Cell culture.** MCF10A immortalized breast epithelial cells<sup>20</sup> were purchased from ATCC and cultured in DMEM/F12 (Invitrogen) supplemented with 10% FBS, 1% pen-strep, 10 ng/mL EGF, 1 µg/mL hydrocortisone, 5 µg/mL insulin, and 100 ng/mL cholera toxin. K562 cells were cultured in RPMI 1640+L-Glutamine (Gibco) supplemented with 10% FBS (Rocky Mountain Biologicals) and 1% pen-strep (Gibco).

**Generating inducible Cas9-expressing MCF10A cell lines.** Lentivirus containing either a doxycycline-inducible or constitutively expressed Cas9 construct were produced by transfecting 293T cells with either pCW-Cas9 (Addgene 50661) or lentiCas9-Blast (Addgene 52962) using the ViraPower Lentiviral Expression System (Thermo) according to manufacturer's instructions. 48 h post-transfection, supernatant was collected and debris removed using a 40 µm syringe filter. MCF10A were transduced with viral supernatant for 48 h and selected with 1 µg/mL puromycin (pCW-Cas9) or 10 µg/mL blasticidin (lentiCas9-Blast) for 96 h. For cells expressing a doxycycline-inducible Cas9, single-cell clones of MCF10A-Cas9 cells were generated by dilution, clones were expanded, and Cas9 expression was confirmed by immunoblotting 96 h following addition of doxycycline at 1 µg/mL. lentiCas9-Blast cells were maintained as a polyclonal line.

pCW-Cas9 cells were used for initial arrayed and pooled screens as well as quantification of editing rates in pHAGE-scKO vector. lentiCas9-Blast cells were used for all CROP-seq experiments.

**Initial tagged transcript cloning method.** Because of high rates of barcode-sgRNA swapping when using this design, we do not recommend use of this protocol.

LentiGuide-puro (Addgene 52963) was modified to confer blasticidin resistance. Puro and its EF-1A promoter were removed via double digest with NEB SmaI (8 h at 25 °C) and MLU1-HF (8 h 25 °C). This product was gel purified using QiaQuick Gel Extraction kit (Qiagen). EF-1A promoter and blasticidin, each with 20 bp homology on both ends, were prepared via PCR from lentiCas9-Blast and gel purified. Fragments were assembled into digested lentiGuide-puro vector using the NEBuilder HiFi DNA Assembly kit with inserts in two-fold molar excess and transformed into NEB C3040H *E. coli* and allowed to incubate overnight at 30 °C. Clones were picked from plate, allowed to grow in LB + amp overnight at 30 °C, and purified using Qiagen Miniprep kit. Individual clones were validated via Sanger sequencing.

Lentiguide-blast was linearized using a digest with BsmB1 (Thermo) at 37 °C for 5 h followed by digestion with SalI HF (NEB) overnight and gel purification. Oligos containing guide sequences and their corresponding barcodes were designed according to the following:

```
tGTGGAAAGGACGAAACACC[G][guide]gttttagagctaG
AAAtagcagagacgCGTCTCAgatcccttggccgcctcccgcg[bar
code]tcgacttaagaccaatgacttaca
```

Where [guide] is a 20 bp guide sequence and [barcode] is an 8 bp barcode sequence uniquely paired to an sgRNA. The [G] included before guide is required for expression from Pol III promoters. Guides/barcodes that generate an extra BsmB1 restriction site when used in this design were excluded. RUNX1 only included four guides because of this filter.

A library of these oligos was ordered as Ultramers from IDT. All oligos were resuspended in water, pooled at equimolar concentrations, and amplified using a 50 µl KAPA HiFi HotStart Ready Mix PCR reaction with 1ng of input DNA. The resulting product was cleaned with a Zymo DNA Clean and Concentrator kit. The purified inserts were assembled into linearized lentiGuide-blast using the NEBuilder HiFi DNA Assembly kit and a molar excess of 1:5 vector to insert. Assembled products were transformed into NEB C3040H *E. coli* and grown overnight at 30 °C in LB + amp. Product was prepared using a plasmid Miniprep kit (Qiagen).

To prepare the insert for the final reaction, a region from the backbone sequence for the CRISPR sgRNA to a region toward the end of the WPRE element was amplified using the KAPA HiFi Hotstart Master Mix and purified using the Zymo Clean and Concentrator kit. The primers used in this reaction add BsmB1 cut sites that generate complementary ends in the final cloning step following digestion. This amplified fragment was ligated into PGEM-T using the PGEM-T kit a clone selected and validation of individual clones by Sanger sequencing. The validated construct was digested with BsmB1 (Thermo) and gel purified.

The fragment isolated from PGEM-T was then ligated into the linearized vector using a 3:1 molar excess of insert to vector using T4 DNA Ligase (New England Biolabs) and overnight incubation at 16 °C. Ligation products were transformed into NEB C3040H (stable) competent cells and grown overnight at 30 °C in LB + amp. Plasmids were recovered using a Plasmid Miniprep kit (Qiagen).

**pHAGE and CROP-seq vector cloning.** The pHAGE\_dsRed\_IRES\_zsGreen vector was modified to contain a multiple cloning site as described in "Quantification of template switching in lentivirus packaging using FACS." The U6-sgRNA cassette containing a 500 bp filler removable by BsmB1 digest was ordered as an IDT gblock. Using the multiple-cloning site, the U6-sgRNA cassette was added in the 3' UTR of the zsGreen/dsRed transgene via Gibson assembly. This vector was modified to remove the zsGreen/IRES/dsRed cassette and replace the CMV promoter with an EF1a promoter.

To clone libraries for this vector or CROP-seq vector (Addgene 86708), the starting vector was digested following the protocol outlined in ref. 21. Oligos corresponding to individual guides with homology for Gibson assembly were ordered as standard DNA oligos from IDT with the following design:

```
[GCCTTATTTTAACTTGCTATTTCTAGCTCTAAAAC][GUIDERC][C][GGTGTTCGTCCTTCCACAAGAT]
```

GUIDERC refers to the reverse complement of the guide sequence. The entire construct may also be reverse complemented, allowing the guide sequence itself to be used rather than the reverse complement.

All oligos were resuspended in water, pooled at equimolar concentrations, and amplified using a 50 µl KAPA HiFi HotStart Ready Mix PCR reaction with 1 ng of input DNA. The following primers were used for amplification:

```
Forward: 5-GCCTTATTTTAACTTGCTATTTCTAGCT-3
Reverse: 5-ATCTTGTGGAAAGGACGAAACA-3
```

These reactions were cleaned with a Zymo DNA Clean and Concentrator kit and cloned into the BsmB1-digested pHAGE vector backbone using the Clontech Infusion HD Cloning Kit.

Ligations were performed using 10 fmol of vector and 200 fmol of double-stranded oligo (1:20 molar ratio of vector to insert). Ligation products were transformed into NEB C3040H (stable) cells according to manufacturer recommendations. Transformations were diluted with 250  $\mu$ L of LB and spread onto 6 LB-AMP plates and incubated at 30 °C for 24 h. Colonies were then scraped into LB, bacterial pellet was collected, and plasmids were recovered using a Plasmid Midiprep kit (Qiagen).

The CROP-seq vector with optimized backbone (CROP-seq-opti) was cloned in a manner similar to the standard CROP-seq vector but with different homology.

Oligos were ordered with the following 3' homology:

5-gtttAagagctaTGCTGGAAACAGCAtagcaagt-3

If ordering in the same format as above (where the oligo is the reverse complement), GCCTTATTTTAACTTGCTATTTCTAGCTCTAAAAC, would be replaced by the reverse complement of the above sequence and amplified with primers: Forward: 5-atcttGTGGAAAGGACGAAACA-3

Reverse: 5-acttgctaTGCTGTTTCCAGC-3

Each of these vectors is also compatible with alternative cloning protocols for lentiGuide-Puro vectors (as long as any homology is adjusted as needed).

#### Quantification of template switching in lentivirus packaging using FACS.

A multiple-cloning site was cloned into pHAGE<sub>dsRed</sub>\_IRES<sub>zsGreen</sub> lentiviral vector between the WPRE and 3 LTR. The multiple-cloning site was assembled from annealing and extension of WPRE<sub>MCS</sub>\_insert<sub>W</sub> and WPRE<sub>MCS</sub>\_insert<sub>R</sub>:

WPRE<sub>MCS</sub>\_insert<sub>W</sub>:

5-ctttggcgcctccccgctggcgccATAACAgctagcTGATGGctcgagcc-3

WPRE<sub>MCS</sub>\_insert<sub>R</sub>:

5-cagctgcctgtaagtcattggtcttaaggtcagCCATCAgctagcTGTTATgg-3

The plasmid was amplified by inverse PCR with pHAGE<sub>WPRE</sub>\_MCS\_GIBS<sub>F</sub> and R:

pHAGE<sub>WPRE</sub>\_MCS\_GIBS<sub>F</sub>

5-TGGctcgagccttaagaccaatgactacaagcagctg-3

pHAGE<sub>WPRE</sub>\_MCS\_GIBS<sub>R</sub>

5-ctagcTGTTATggcgccccagcggggagggcgcccaaag-3

The fragments were cloned by Gibson assembly. Clones of pHAGE<sub>dsRed</sub>\_IRES<sub>zsGreen</sub>\_WPRE<sub>MCS</sub> were chosen by Sanger sequencing and expression of the fluorescent proteins after transfection and lentiviral packaging.

To make pHAGE EBFP or EGFP<sub>IRES</sub>\_dsRed<sub>WPRE</sub>\_MCS, pHAGE<sub>dsRed</sub>\_IRES<sub>zsGreen</sub>\_WPRE<sub>MCS</sub> was cut with BamHI and ClaI to remove the zsGreen and IRES. The ends were blunted and religated to make pHAGE<sub>dsRed</sub>\_WPRE<sub>MCS</sub>. EGFP or EBFP (amplified with eGFP<sub>gibsF</sub> and eGFP<sub>IRES</sub>\_GibsR) and an IRES (IRES<sub>GibsF</sub>, IRES<sub>GibsR</sub>) were cloned into the NotI site 5' of the dsRed by Gibson assembly. EBFP was ordered as a gblock from IDT with 3 nt changes from EGFP. Correct clones were identified by sequencing. The dsRed is not expressed in this construct.

eGFP<sub>gibsF</sub>:

5-gccatccagctgttttgactccatagaagaccggcATGGTGAGCAAGGGCGAGGAG-3

eGFP<sub>IRES</sub>\_GibsR:

5-ggatccCTACTTGTACAGCTCGTCCATGCCG-3

IRES<sub>GibsF</sub>:

5-ATCACTCTCGGCATGGACGAGCTGTACAAGTAGgg  
atccctccccccccctaacgttac-3

IRES<sub>GibsR</sub>:

5-ctccttgatgacgtctcggaggaggccatggcgccatgtgtggccattatcatcgtgttttcaagg-3

EBFP

5-ATGGTGAGCAAGGGCGAGGAGCTGTTACACCGGGTGGTGCCCATCCTGGTTCGAGCTGGACGGCGACGTAAACGCCACAAGTTCAGCGTGTCCGGCGAGGGCGAGGGCGATGCCACCTACGGCAAGCTGACCCTGAAGTTCATCTGCACCACCGCAAGCTGCCCGTGCCTGGCCACCCTCGTGACCACCTGACCCACGGCGTGCAGTGTTCAGCCGCTACCCCGACCACATGAAGCAGCAGCACTTCTTCAAGTCCGCCATGCCCGAAGGCTACGTCCAGGAGCGCACCATCTTCTTCAAAGGACGCAACTACAAGACCCGCGCCGAGGTGAAGTTCGAGGGCGACACCCTGGTGAACCGCATCGAGCTGAAGGGCATCGACTTCAAGGAGGACGGCAACATCCTGGGGCAAGAGTGGAGTACAACCTTAAACAGCCACAACGTCTATATCATGGCCGACAAGCAGAAGAACGGCATCAAGGTGAACCTCAAGATCCGCCACAACATCGAGGACGGCAGCGTGCAGCTGCCGACCACCTACCAGCAGAACACCCCATCGGCGACGGCCCGTGTCTGCTGCCGACAACCACTACCTGAGCACCCAGTCCGCCCTGAGCAAAGACCCCAACGAGAAGCGCGATCACATGGTCTGCTGGAGTTCGTGACCGCCGCCGGGATCCTCTCGGCATGGACGAGCTGTACAAG-3

15 bp barcodes (lenti-barcode and lenti-barcode-r) were cloned into the multiple-cloning site between the WPRE and 3' LTR for both the EBFP and EGFP constructs by Gibson assembly. Single clones were prepared and the barcode identified by Sanger sequencing.

lenti-barcode:

5-atctccctttggcgcctccccgctgggGGATCCAGNNNNNNNNNNNNNNNtgcgaccttaagaccaatgactacaagg-3

lenti-barcode-r:

5- CCTTGTAAGTCATTGGTCTTAAAGGCTCGA -3

Lentivirus was packaged by transfection of barcoded EGFP or EBFP constructs either alone or in an equimolar mix along with helper plasmids (pHDM-Hgpm2, pHDM-Tatlb, pRC-CMVRev1b, and pHDM-VSV-G) into HEK293T cells using Lipofectamine 2000 (Invitrogen). Viral supernatant was collected after 48 h, spun to remove debris, snap frozen in liquid nitrogen, and stored at -80 °C. To titer the packaged lentiviruses, they were thawed on ice and added to MCF10A cells with media containing 8  $\mu$ g/ml polybrene, and the frequency of transduced cells 48 h post-transduction was determined by flow cytometry.

To sort blue+ and green+ populations, 400,000 of MCF10A TP53 cells (Horizon Discovery) in 5 ml media plus 8  $\mu$ g/ml polybrene were transduced at a MOI ~0.1, with either of the EGFP or EBFP expressing viruses that had been packaged singly, a mix of the EGFP and EBFP expressing viruses that had been packaged singly, or the EGFP and EBFP expressing viruses that had been packaged together. The cells were cultured for 4 weeks to avoid residual plasmid contamination following transduction. An equal number of cells transduced with EGFP and EBFP virus were mixed to determine the rate of contamination resulting from FACS error. The mixed cells along with others were sorted for blue+ or green+ populations using a FACS Aria II (Becton Dickinson) that had been compensated for the overlap

between the EBFP and EGFP emission spectra. Genomic DNA was harvested from each population using the Qiagen DNeasy kit, and barcodes were amplified from 2–36 ng of genomic DNA in 50  $\mu$ l Robust polymerase (Kapa) reactions with primers `bwds_p5_WPRE_BC_F` and `bwds_next_WPRE_BC_R`.

```
bwds_next_WPRE_BC_R:  
GGCTCGGAGATGTGTATAAGAGACAG  
5-gaaatcatcgtcctttccttgct-3  
bwds_p5_WPRE_BC_F:  
5-AATGATACGGCGACCACCGAGAgcgccgatgcttgaagtc  
attggtcttaaggctc-3
```

PCR products were purified with Ampure (Agilent) and P7 index sequences added by an additional six cycles of PCR. PCR products were purified, quantified, pooled, and single-end sequenced on an Illumina Nextseq500 with Read1 primer `bwds_WPRE_bc_seqF` and standard Illumina i7 primers.

```
bwds_WPRE_bc_seqF:  
5-GCGCCGATGCCTTGTAAGTCATTGGTCTTAAAGG  
CTCGA-3
```

**Analysis of FACS data from pHAGE-GFP and pHAGE-BFP experiments.** Background percentage of contaminating barcodes in the BFP/GFP-sorted cells from the mixed cells control was subtracted from numbers obtained for the pooled virus samples. Fraction of GFP cells, determined from FACS gating, was fixed; and the expected fraction of barcode contamination in the BFP and GFP was simulated. Note that the expected contamination of green barcodes in the BFP-sorted cells is the template-switching rate multiplied by the fraction of green cells. The expected rate of contamination of BFP barcodes in the GFP-sorted cells is the template-switching rate multiplied by the BFP fraction (1 – GFP fraction). Sum of the squared error between observed and expected values for rates of contamination was calculated for a range of different lentivirus swap rates, and minimal value was taken to be the most likely swap rate.

Note that, unlike a library of plasmids, in a mix of two plasmids, only half of all chimeric products will be detectable as many virions will be homozygous (i.e., contain the same construct, and thus chimeric products are identical to the original). To give an analogous example, in a barnyard experiment for a single-cell assay, mouse–mouse or human–human multipliers cannot be detected and thus estimated rates of ‘doublets’ have to be adjusted accordingly. When the plasmids are equimolar and the swap rate is 50%, for example, one would expect to observe a 75% rate of the intended barcode and a 25% rate of the unintended barcode. This ratio will change according to the molar concentration of the two plasmids. In **Figure 1e**, we assume the pool was composed of 61.7% GFP plasmid, corresponding to the fraction of GFP<sup>+</sup> cells relative to the total number of GFP<sup>+</sup> and BFP<sup>+</sup> cells  $4.59/(4.59 + 2.85)$  or 61.7% as explained in **Supplementary Figure 5**. This analysis was also performed without fixing the fraction of GFP<sup>+</sup> cells to the value measured by FACS to ensure results were concordant (**Supplementary Fig. 6**). The minimum sum of squared error over the grid of simulated lentivirus swap rate and fraction of GFP cells were taken to be the most likely set of parameter values.

**CRISPRi experiment.** K562 expressing dCas9-BFP-KRAB (gift of the Bassik lab, Addgene 46911) and MCF10A expressing dCas9-BFP-KRAB (made by transduction with lenti\_UCOE\_EF1-dCas9-BFP-KRAB, plasmid, a gift of the Weissman lab (available

on Addgene soon; see <https://weissmanlab.ucsf.edu/CRISPR/CRISPRiacellineprimer.pdf>) were transduced with lenti-mCherry under control of a CAG promoter (pCAG\_mCherry pKH143, gift of the Bassik lab, unpublished), and sorted such that the resulting population is enriched for mCherry expression.

A spacer targeting the CAG promoter was cloned into the KHH030 (Addgene 89358), CROP-seq, and pHAGE-scKO sgRNA expression vectors. The CROP-seq and pHAGE-scKO vectors were modified by Q5-Site Directed Mutagenesis (NEB) to use the previously described sgRNA-(F+E)-combined optimized backbone<sup>22</sup> (we refer to this as CROP-seq-opti). The CRISPRi mCherry<sup>+</sup> K562 and MCF10A cells were transduced with the CAG-targeting sgRNA and assayed for mCherry.

All viruses for the CRISPRi experiments were made by the Co-operative Center for Excellence in Hematology Vector Production core. All sorting was performed on a FACS Aria II (Becton Dickinson).

**Editing-rate experiment for pHAGE-scKO.** To confirm that our pHAGE-scKO vector exhibited reduced editing efficiency, we performed editing with a guide to *TP53* from our screen (GAGCGCTGCTCAGATAGCGA) in both lentiGuide-Blast and pHAGE-scKO using our pCW-Cas9 MCF10A cells. Cells were passaged for 18 d after induction of Cas9 expression with dox, and gDNA was harvested using Qiagen DNeasy kit and amplified using primers CTAAATGGCTGTGAGAGAGCTCAGCCACACGCAAATTTCTTCC and ACTTTATCAATCTCGCTCCAAA CCCCCTGCCCTCAACAAGATGT. These were then amplified using KAPA HiFi Hotstart Ready Mix (KAPA) using the following indexed primers: AATGATACGGCGACCACCGAGATCTACA CagctaggcCTAAATGGCTGTGAGAGAGCTCAG

```
CAAGCAGAAGACGGCATAAGAGAT[INDEX]gacctcggcA  
CTTTATCAATCTCGCTCCAAACC
```

Libraries were sequenced on MiSeq, and reads were then processed using the method described in McKenna *et al.*<sup>23</sup>.

Briefly, reads are trimmed of low-quality bases using Trimmomatic, merged using Flash, aligned to the reference of the locus surrounding the guide using needle, and unique genotypes are quantified. The wild-type genotype fraction was taken to be the proportion of non-wild-type alleles. We did not use UMIs in this experiment, and thus it may overestimate editing rate.

**Knockout experiments.** For all screens, each plasmid library was transfected along with plasmids provided with the ViraPower Lentiviral Expression into 293T cells. At 48 and 72 h post-transfection, supernatant was collected and filtered using a 40  $\mu$ m steriflip filtration system (EMD Millipore). For arrayed experiments, individual plasmids were transfected and viruses produced as described above. For pHAGE-scKO and arrayed/pooled pLGB-scKO vector experiments, virus was concentrated using Peg-it virus concentration solution (SBI). Viral titer of the concentrated lentiviral library was determined by transduction of MCF10A-Cas9 cells for 48 h at several viral dilutions, splitting cells into replica plates, and subjecting replica plate to blasticidin. Percent control growth was used to assess MOI. MCF10A-Cas9 cells with estimated MOIs of 0.3 carried forward.

For pHAGE-scKO and arrayed/pooled pLGB-scKO vector experiments, media were switched to 1  $\mu$ g/mL doxycycline to induce expression of Cas9 in pCW-Cas9 cells. LentiCas9-Blast

cells were used for CROP-seq experiments. Editing was allowed take place for 14 d for arrayed and pooled pLGB-scKO and 21 d for pHAGE-scKO and CROP-seq experiments. Media were changed every 48 h, and cells were cultured every 96 h. For the first half of editing, cells were cultured in the presence of 5 µg/mL blasticidin and 0.5 µg/mL puromycin to ensure high sgRNA and Cas9 expression. In all CROP-seq KO experiments (but not our CRISPRi experiment), we used the CROP-seq vector from Datlinger *et al.*<sup>7</sup> without modification (Addgene 86708).

**Doxorubicin treatment.** After editing, MCF10a cells were seeded in 10 cm plates at  $1 \times 10^6$  cells per well, allowed to attach overnight, and media replaced with MCF10A media alone (mock) or MCF10A media containing 500 (arrayed and pooled pLGB-scKO experiments) or 100 nM (pHAGE-scKO and CROP-seq experiments) doxorubicin prepared from a 500 µM stock of doxorubicin (Sigma) in water. 24 h after drug exposure, untreated and doxorubicin-treated cells were harvested by trypsinization, washed with PBS, and used for downstream assays.

**Single-cell RNA sequencing.** Cells were captured using one lane of a 10X Chromium device per sample using 10X V1 Single Cell 3'-Solution reagents (10X Genomics). Approximately 4,000–7,000 cells were captured per lane for each condition. Protocols were performed according to protocol, holding 10–30 ng of full-length cDNA out of downstream shearing and library prep steps in order to provide material for barcode-enrichment PCR.

Final libraries were sequenced on NextSeq500. 10X V1 samples were sequenced using the following read configuration:

R1: 64, R2: 5, I1: 14, I2: 8

Our initial arrayed and pooled doxorubicin-treated samples using pLGB-scKO were aggregated using cellranger aggregate to normalize the average number of mapped reads per cell. This yields an average of 37,732 reads per cell; 2,263 median genes per cell; and a median of 8,279 UMIs per cell.

Our CROP-seq mock sample was sequenced to an average depth of 120,797 raw reads per cell in 6,598 cells. A median of 4,619 genes per cell was detected and a median UMI count of 22,495 per cell. Our CROP-seq doxorubicin-treated sample was sequenced to an average depth of 123,445 raw reads per cell in 6,283 cells. A median of 3,500 genes per cell was detected, and we observed a median UMI count of 15,324 per cell. At this depth the average duplication rate is approximately 78%.

**Enrichment PCR.** For all experiments, a heminested PCR starting from 5 ng of full-length cDNA was used to enrich for barcodes that assign a target to each cell. PCR reactions were performed with a P7 reverse primer (as introduced by the 10X Chromium V1 oligo DT RT primer). Importantly, the protocol for the 10X V2 protocol (not used here) would be different—see <https://github.com/shendurelab/single-cell-ko-screens#enrichment-pcr> for more information. For pHAGE-scKO and pLGB-scKO, the first PCR was performed with:

5-TCCTGGGATCAAAGCCATAGT-3

and for CROP-seq:

5-TTTCCCATGATTCTTCATATTTGC-3

as the forward primer, priming to the blasticidin transcript with no nontemplated sequence for pHAGE-scKO and pLGB-scKO, and to part of the U6 promoter in CROP-seq. For pLGB-scKO the second PCR was performed with:

5-TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGG  
ACGAGTCGGATCTCCCTT-3

for pHAGE-scKO with:

5-TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGA  
ACGGACTAGCCTTATTTTAACTTG-3

and for CROP-seq with:

5-TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGcT  
TGTGGAAAGGACGAAACAC-3

as the forward primer, priming adjacent to the barcode/guide sequence in each design and adding the standard Nextera R1 primer. Samples were indexed in a final PCR using standard Nextera P5 index primers of the form:

5-AATGATACGGCGACCACCGAGATCTACAC[8bp  
Index]TCGTCGGCAGCGTC-3

PCRs were cleaned with a 1.0X AmpureXP cleanup and 1 µl of a 1:5 dilution of the first PCR and 1:25 dilution of the second PCR were carried in each reaction.

**Digital gene expression quantification.** Sequencing data from each sample was processed using cellranger 1.3.1. Each lane of cells was processed independently using cellranger count, aggregating data from multiple sequencing runs. For the comparison between arrayed and pooled screens, cellranger aggregate was used to downsample data from each screen to an equal average number of mapped reads.

**Assigning cell genotypes.** Barcode-enrichment libraries were separately indexed and sequenced as spike-ins alongside the whole-transcriptome scRNA-seq libraries. Final UMI and cell-barcode assignments were made for each read by processing these samples with cellranger 1.3.1, as was done for the whole-transcriptome libraries.

A whitelist of guide or target barcode sequences was constructed using all guides or target barcodes in the library. For each read in the position-sorted BAM file output by cellranger 1.3.1, the final cell barcode and UMI are extracted. If either of these fields is not populated, indicating a problem with the sequence, the read is ignored. Using the cDNA read, we attempt to find a perfect match for the sequence preceding the guide or barcode (GTGGAAAGGACGAAACACCG for CROP-seq and CGCCTCCCGCG for pLGB-scKO). If a perfect match is not found, we attempt to locate the sequence using a striped Smith–Waterman alignment. If a match or alignment is found, the guide or barcode sequence is extracted. If the extracted sequence does not perfectly match a whitelist sequence, we search for a matching whitelist sequence within an edit distance of half the minimum edit distance between any pair of guides or barcodes in the library (rounded down). If no match is found, the molecule is ultimately discarded. Matches to the whitelist are tracked for each cell.

We also remove likely chimeric sequences using the approach outlined in Dixit<sup>24</sup>. Briefly, within each cell we calculate the number of times a given UMI is observed with each observed guide assignment. We then divide these counts by the total instances of the respective UMI across all observed guide assignments within that cell. For UMI-guide assignment combinations where this fraction is less than 20%, we do not count the UMI toward the final observed guide assignment counts. While this has some impact on the raw data, we find the benefits to be modest.

To make a set of final assignments, we take all whitelist sequences that have over ten reads and account for over 7.5%



of the whitelist reads assigned to a given cell, where multiple sequences can be assigned to each cell. This set of assignments is merged with the filtered gene expression matrices output by cellranger such that only assignments to the filtered cells appear in the final data set.

Note that when processing CROP-seq data without PCR enrichment, we lowered the requirement for reads supporting a given guide to 3 to account for the decreased coverage of these transcripts.

**Estimation of multiplicity of infection and capture rate.** The most likely multiplicity of infection (MOI) and capture rate given the distribution of guide counts per cell were estimated using the generative model described in ref. 5. Briefly, a log likelihood is calculated using a zero-truncated poisson (MOI postselection) convolved with a binomial (incomplete capture of barcoded transcripts). This model is used to estimate the most likely set of MOI and capture rate values.

**Monocle2 usage.** PCA + tSNE, density peak clustering, differential expression testing, and size-factor estimation were performed using the monocle2 (ref. 25) functions `reduceDimension`, `clusterCells`, `differentialGeneTest`, and `estimateSizeFactors` unless otherwise noted.

**Removing low-quality cells.** We consistently observed a cluster of cells with much lower UMI counts on average than the rest of the data set when performing dimensionality reduction. To avoid including these cells in downstream analysis, we perform a simple procedure to remove any cluster with low average UMI counts. We perform PCA followed by TSNE on genes expressed in at least 50 cells for each condition, perform density peak clustering on two-dimensional tSNE space, calculate the average size factor over each cluster, and filter out clusters of cells with an average size factor of  $2^{-0.85}$  or lower before downstream analysis.

**Simulating loss in power from barcode swapping.** Assignments were permuted for a fraction of cells ranging from 0 to 100% and kept fixed for the remaining fraction of cells. We tested for genes differentially expressed across the target assigned to each cell (testing genes detectably expressed in at least 50 cells; full model  $\sim$ target\_gene). Differentially expressed genes at FDR of 5% were counted. Ten samplings were performed for each swap rate.

For the simulation performed on our own data, cells with a single target assignment from 100 nM doxorubicin-treated cells in our CROP-seq experiment were taken as the starting set of cells.

For the simulation on data from Adamson *et al.*<sup>4</sup>, processed data were obtained from GEO (GSE90546). Assignments of cells to targets were used as provided on GEO, and only cells noted as having high-quality assignment to a single target were used. Because of the large number of cells (50,000+) in the UPR experiment from this study and the large number of differential tests required for these simulations, the number of cells assigned to each target was downsampled two-fold to reduce runtime. We also performed tests on a data set further downsampled to approximately 6,000 cells to illustrate the impact of initial power.

**tSNE embedding demonstrating TP53-enriched cluster.** 20 dimensions from PCA were carried into tSNE to two dimensions.

All cells, including cells with guides to multiple targets and no assigned target, were included in dimensionality reduction for this plot. Percentages of cells with guides to *TP53* and *ARID1B* were calculated, including cells that contain guides to multiple targets. All cells with *TP53* guides were counted as *TP53* cells only.

#### **Enrichment of tumor suppressors in specific molecular states.**

Only cells containing a guide to a single target were considered in enrichment testing. A Chi-squared test was used to determine whether the distribution of individual sgRNAs and targets in tSNE space was significantly different from nontargeting controls at 5% FDR. Targets which did not pass this test and did not have an individual sgRNA pass the test were excluded from the subsequent enrichment tests. For each sgRNA of the remaining targets, we sought to estimate the functional editing rate (probability of a cell having a true LoF given that it received that sgRNA). Such estimates would be confounded if accounting for the possibility of edits that cause LoF for the target gene but have incomplete penetrance on the cellular phenotype. Therefore, we used an expectation-maximization approach to estimate the functional edit rate of each sgRNA relative to the unknown functional edit rate of the most efficient sgRNA for a given target.

The t-SNE cluster distribution of all cells in which a given sgRNA was detected was modeled as a mixture of the t-SNE cluster distribution of cells with a functional edit for the sgRNA's target gene and the t-SNE cluster distribution of nontargeting controls, where the mixing parameter is the relative functional edit rate for that sgRNA. In the expectation step, the t-SNE cluster distribution of cells with a functional edit for the target is estimated as the weighted average of the empirical t-SNE cluster distributions of each sgRNA for the target, weighted by the current estimates of the relative functional edit rate of the sgRNAs. In the maximization step, the relative functional edit rate of each sgRNA for the target is chosen to maximize the likelihood of the observed t-SNE cluster distribution for cells receiving that sgRNA under the multinomial mixture model.

After estimating the relative functional edit rate for each sgRNA, a weighted contingency table was constructed where the rows are targets, the columns are t-SNE clusters, and the values are weighted cell counts, and where a cell's weight is proportional to the relative functional edit rate for the sgRNA it received. Fractional values were rounded down. Fisher's exact test was applied to this weighted contingency table to test for enrichment of targets amongst t-SNE clusters. Targets were defined as enriched at an FDR of 10%. Chi-square and Fisher's exact test were performed using R functions `chisq.test` and `fisher.test`, respectively.

#### **Principal component and gene set enrichment analysis.**

Pairwise differential gene expression analysis was performed between enriched target cells and nontargeting controls for cells in all significantly enriched target-cluster pairs from our enrichment testing. The union of all differentially expressed genes across targets (FDR 5%) was used to perform principal component analysis. Gene set enrichment analysis was performed on genes that had the highest positive and negative loadings for principal component 1 (less than  $-0.02$  or greater than  $0.02$ ). Gene set enrichment analysis was performed using the piano R package and the hallmarks gene set from MSigDB. Gene sets were

defined as enriched at an FDR of 1%. PCA was performed using the `prcomp` function in R.

**Code availability.** Code and information on how to access additional data files relevant for secondary analysis can be found on Github at <https://github.com/shendurelab/single-cell-ko-screens>.

**Life Sciences Reporting Summary.** Further information on experimental design is available in the **Life Sciences Reporting Summary**.

**Data availability.** Data is available on GEO via accession [GSE108699](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE108699) and is also provided via the Github repository described in “Code availability.” pHAGE-GFP, pHAGE-BFP, and the CROP-seq vector with the CRISPRi-optimized backbone sequence described in the Online Methods are available

on Addgene as [106281](https://www.addgene.org/106281), [106282](https://www.addgene.org/106282), and [106280](https://www.addgene.org/106280). All CROP-seq experiments, except for the one presented in **Figure 2d**, were carried out with the original CROP-seq vector described in ref. 7. For the experiments shown in **Figure 2d**, we used our own version of CROP-seq modified to contain a backbone optimized for CRISPRi, available on Addgene as described above.

20. Debnath, J., Muthuswamy, S.K. & Brugge, J.S. *Methods* **30**, 256–268 (2003).
21. Sanjana, N.E., Shalem, O. & Zhang, F. *Nat. Methods* **11**, 783–784 (2014).
22. Chen, B. *et al. Cell* **155**, 1479–1491 (2013).
23. McKenna, A. *et al. Science* **353**, aaf7907 (2016).
24. Dixit, A. Preprint at <https://www.biorxiv.org/content/early/2016/12/12/093237> (2016).
25. Qiu, X. *et al. Nat. Methods* **14**, 309–315 (2017).

## Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form is intended for publication with all accepted life science papers and provides structure for consistency and transparency in reporting. Every life science submission will use this form; some list items might not apply to an individual manuscript, but all fields must be completed for clarity.

For further information on the points included in this form, see [Reporting Life Sciences Research](#). For further information on Nature Research policies, including our [data availability policy](#), see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

## ▶ Experimental design

## 1. Sample size

Describe how sample size was determined.

The number of cells for single cell RNA-Seq were determined by obtaining a reasonable amount of coverage in terms of minimum number of cells per target (roughly more than 50 cells per genotype).

## 2. Data exclusions

Describe any data exclusions.

No data exclusions

## 3. Replication

Describe whether the experimental findings were reliably reproduced.

Similar signatures were observed across initial arrayed and CROP-Seq pooled experiments. CRISPRi knockdown of mCherry was performed on multiple cells lines with multiple controls. GFP/BFP experiments were done sorting for both BFP and GFP to ensure results were symmetrical.

## 4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

Not applicable

## 5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

Not applicable

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

## 6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or in the Methods section if additional space is needed).

n/a Confirmed

- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.)
- A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- A statement indicating how many times each experiment was replicated
- The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section)
- A description of any assumptions or corrections, such as an adjustment for multiple comparisons
- The test results (e.g.  $P$  values) given as exact values whenever possible and with confidence intervals noted
- A clear description of statistics including central tendency (e.g. median, mean) and variation (e.g. standard deviation, interquartile range)
- Clearly defined error bars

See the web collection on [statistics for biologists](#) for further resources and guidance.

## ► Software

Policy information about [availability of computer code](#)

### 7. Software

Describe the software used to analyze the data in this study.

The single cells analysis packages cellranger (10X Genomics) and Monocle2 were used in this article. Custom analysis software is being prepared for distribution on github and will be available upon request for review.

For manuscripts utilizing custom algorithms or software that are central to the paper but not yet described in the published literature, software must be made available to editors and reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). *Nature Methods* [guidance for providing algorithms and software for publication](#) provides further information on this topic.

## ► Materials and reagents

Policy information about [availability of materials](#)

### 8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

No restrictions.

### 9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

No antibodies were used in the study.

### 10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

MCF10A breast epithelium cells were purchased from ATCC (CRL-10317), MCF10A TP53 -/- cells were purchased from Horizon Discovery (HD 101-005) and K562 cells were a gift from the Bassik lab.

b. Describe the method of cell line authentication used.

MCF10A and MCF10A TP53 -/- cells were not authenticated but used within 10 passages of purchase. K562 cells were not authenticated.

c. Report whether the cell lines were tested for mycoplasma contamination.

MCF10A cell lines were tested and confirmed negative for mycoplasma contamination. K562 and TP53 -/- cells were not tested for mycoplasma contamination.

d. If any of the cell lines used are listed in the database of commonly misidentified cell lines maintained by [ICLAC](#), provide a scientific rationale for their use.

None of the cell lines used in this study are registered in the ICLAC database of commonly misidentified lines.

## ► Animals and human research participants

Policy information about [studies involving animals](#); when reporting animal research, follow the [ARRIVE guidelines](#)

### 11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

No animals were used in the study.

Policy information about [studies involving human research participants](#)

### 12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

No human subjects were used in the study.

## Flow Cytometry Reporting Summary

Form fields will expand as needed. Please do not leave fields blank.

### ▶ Data presentation

For all flow cytometry data, confirm that:

- 1. The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- 2. The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- 3. All plots are contour plots with outliers or pseudocolor plots.
- 4. A numerical value for number of cells or percentage (with statistics) is provided.

### ▶ Methodological details

- |  |   |
|--|---|
| 5. Describe the sample preparation.  | The cells are MCF10A with TP53 deleted from Horizon Discovery. The cells were transduced with indicated lentiviruses and cultured for an additional 4 weeks in the recommended media. For flow cytometry, 6 cm plates of cells were removed from the plate with 0.25% trypsin, washed with PBS and resuspended in PBS + 1% heat-inactivated FBS, 1mM EDTA, 25 mM HEPES pH 7.5.                  |
| 6. Identify the instrument used for data collection.                                   | Becton Dickinson FACS Aria II   |
| 7. Describe the software used to collect and analyze the flow cytometry data.          | The data was collected using FACSDiva version 8 software. Data was analyzed using FlowJo 10   |
| 8. Describe the abundance of the relevant cell populations within post-sort fractions. | Cells numbers and percentages with post-sort fraction are listed in Figure S4.  |
| 9. Describe the gating strategy used.  | Before analysis of fluorescence, live, single cells were gated using FSC-A and SSC-A (for intact cells) and FSC-A and FSC-H (to ensure that only singlets were analyzed). The green+ and blue+ gates were set after compensating for the overlap between the EGFP and EBFP emission using negative and singly positive cells. Those gates were set to exclude non and double-fluorescent cells. |

Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.