

1 CRISPR/Cas9 screening using unique molecular identifiers

2 **Bernhard Schmierer^{1,#}, Sandeep K. Botla^{1,#}, Jilin Zhang¹, Mikko Turunen², Teemu Kivioja² and Jussi**
3 **Taipale^{1,2,*}**

4 ¹Department of Medical Biochemistry and Biophysics, Karolinska Institutet, 171 77 Stockholm, Sweden.

5 ²Genome-Scale Biology Research Program, Faculty of Medicine, University of Helsinki, PO Box 63
6 FI-00014 Helsinki, Finland.

7 * corresponding author

8 # equal contribution

9

10 **Loss of function screening by CRISPR/Cas9 gene knockout with pooled, lentiviral guide libraries is a**
11 **widely applicable method for systematic identification of genes contributing to diverse cellular**
12 **phenotypes. Here, random sequence labels (RSLs) are incorporated into the guide library, which act as**
13 **unique molecular identifiers (UMIs) to allow massively parallel lineage tracing and lineage dropout**
14 **screening. RSLs greatly improve the reproducibility of results by increasing both the precision and the**
15 **accuracy of screens. They reduce the number of cells and sequencing reads needed to reach a set**
16 **statistical power, or allow a more robust screen using the same number of cells.**

17

18 Pooled CRISPR/Cas9 loss of function screening is a powerful approach to identify genes contributing to a
19 wide range of phenotypes¹. A library of guide sequences is integrated lentivirally into Cas9-expressing cells,
20 which are then subjected to a selection pressure. Relative guide frequencies in the population before and
21 after selection are quantified by next generation sequencing (NGS) to determine depleted and enriched
22 guides.

23 The approach has been applied successfully²⁻⁵, but suffers from several shortcomings: First, the
24 presence of a guide does not necessarily cause loss of the corresponding gene, and cells sharing the same
25 guide have distinct genotypes and phenotypes. Second, identification of guides that are under negative
26 selection can be confounded by random drift and undersampling. Third, growth characteristics of
27 individual cells can vary substantially^{6,7} and the site of viral integration can affect the phenotype. For these
28 reasons, each guide needs to be present in a large number of cells. In conventional screens, only the sum
29 of all cells with a specific guide is measured, and no information regarding the distribution of cell behaviors
30 can be obtained. Optimal identification of hit genes would require a method that individually tracks clonal
31 lineages derived from single virus-transduced cells.

32 Here, we address these issues by incorporating an RSL into the guide-library plasmid (**Fig. 1a**) to
33 allow tracing of thousands of individual virus-transduced cell lineages in a CRISPR screen. Depending on
34 the kinetics of editing, these unique molecular identifiers (UMIs)⁸ either trace single clones of identically
35 edited cells, or small populations of sublineages with different editing outcomes at the same locus
36 (**Supplementary Fig. 1a**). Such massively parallel lineage tracing enables both lineage dropout analysis
37 (LDA), and the creation and analysis of internal replicates (IRA), while retaining the option of conventional,
38 total read count analysis (TCA, **Fig. 1b**). Perhaps counterintuitively, analysis of hundreds to thousands of

39 RSL-labelled cell lineages per guide neither requires more cells per guide, nor deeper sequencing. The RSL
40 approach simply splits the total guide read count obtained to read counts representing individual
41 constituent cell lineages, thus increasing the amount of information that is obtained, and consequently
42 improving both precision and accuracy of the screen.

43 To demonstrate the power and flexibility of the approach, we screened the human colorectal
44 carcinoma cell line RKO for essential genes with an RSL-guide library targeting 2325 genes with 10 guides
45 per gene⁵. For experimental details see **Supplementary Information**. Briefly, Cas9 expressing RKO cells
46 were transduced with the guide library, and samples were taken at Day 4 and Day 28 after transduction
47 (control and treatment time points, respectively). Guide frequencies in the two samples were then
48 assessed by NGS. The experiments were run at far larger screen size (we define “screen size” as the number
49 of cells per guide sequence) and sequencing depth (reads per guide) than previous screens^{4,5}. Such
50 redundancy allows subsequent subsampling using RSL information, and robust testing of different
51 analytical methods at varying screen sizes (**Fig. 1c**).

52 The plasmid library input contained 78 million unique RSL-guide combinations, 93% of which were
53 also detected in the virus-transduced cell populations (**Supplementary Fig. 1b**). Based on the Poisson-
54 distribution, this indicates that about half of the RSL-guides were incorporated into one or two cell
55 lineages. Because only a subset of the cells can be harvested at each time-point, undersampling is
56 unavoidable, and some RSL-guides and thus cell lineages were present only in one of the time points
57 (**Supplementary Fig. 1b, right**). Such undersampling and loss of cell lineages occurs whether or not RSLs
58 are present, however goes undetected in their absence. With RSLs, the effect becomes apparent and can
59 be used in quality control of individual experiments as well as in filtering out inconsistently sampled
60 lineages prior to data analysis.

61 RSL-labelled, distinguishable guide sequences can be used to split the data into internal replicates,
62 which in turn allow the usage of classical statistical tools to test for significant differences. To demonstrate
63 the approach, RSL-guides were binned into 64 internal replicates per guide. The median effect size (**Fig.**
64 **1d**) as well as a median-based version of strictly standardized mean difference (SSMD)⁹ were then used to
65 rank the guides (internal replicate analysis using SSMD, IRA/SSMD, **Supplementary Fig. 1c**). In addition,
66 RSL-labelled guides enable lineage dropout screening, where gene hits are called solely based on the
67 number of lost RSL-guide lineages (lineage dropout analysis, LDA, **Fig. 1e**).

68 To evaluate IRA/SSMD and LDA, and to compare them with conventional TCA performed with the
69 pipeline MAGeCK¹⁰, we assessed the ranks of a set of known essential genes (accuracy), and the hit gene
70 overlap between experimental replicates (precision). In principle, RSL-based methods should outperform
71 TCA when the number of cells per guide is relatively low, and their benefit should progressively decrease
72 as the number of cells per guide approaches infinity. Thus, the comparisons were performed using the
73 complete dataset, and subsamples of the data that were similar in sample size to published screens^{4,5}
74 (**Fig. 2**).

75 Both IRA/SSMD and LDA were more accurate than TCA, as indicated by lower hit ranks of 20 known
76 essential, ribosomal proteins (**Fig. 2a, Supplementary Fig. 1d**). Both IRA/SSMD and LDA were also more
77 precise than TCA, with much improved replicate concordance between the top-ranked 5% of genes (**Fig.**
78 **2b**). Consistently with the theoretical considerations, our analysis revealed that the RSL-based methods

79 were far more robust at smaller screen sizes than TCA. Even when the screen size was approximately one
80 order of magnitude larger, TCA was still inferior to either RSL-based method. At smaller screen size, the
81 number of highly significant hit genes (FDR < 1%) was massively increased in lineage dropout analysis when
82 compared to total readcount analysis (**Fig. 2c**).

83 To summarize, RSLs dramatically improve accuracy, precision, and statistical power in CRISPR/Cas9
84 screening. The RSL strategy is not limited to CRISPR knockout screening, but can be applied in other
85 screening methods such as CRISPR-dependent inhibition or activation screens^{2, 11}. We expect the RSL
86 method to become instrumental in the interrogation of small genomic features, e.g. exons, promoters,
87 and even individual transcription factor binding sites. In many of these cases there is just one possible
88 guide sequence, and the inclusion of RSLs is the only way to obtain the replicates that are required for hit
89 calling. In the absence of precise knowledge of both on- and off-target activity, inclusion of multiple guide
90 positions is however still important, and rescue experiments and/or analysis of the mutational spectrum
91 of the cutsite are necessary to establish that the mutation induced by the guide results in the observed
92 phenotype. Incorporation of RSLs is technically straightforward, and does not require a higher number of
93 cells or sequencing reads compared to conventional approaches. In contrast, RSLs give the same statistical
94 power at lower number of cells and/or sequence reads, improving the economy of CRISPR/Cas9 screens.
95 Conversely, RSLs improve accuracy and precision at a given number of cells per guide, which is particularly
96 advantageous in cases where cell numbers are limiting, such as in primary cells.

97

98 **Figure Legends.**

99

100 **Figure 1.**

101

102 **a. Library design. Top.** Guide plasmid. The plasmid library contains an i7 index read primer and an
103 inert, untranscribed RSL downstream of the tracrRNA. **Bottom.** Sequencing library. Sequencing
104 was done with a custom primer (Seq) placed directly upstream of the guide (gRNA). The sample
105 index and RSL were read as two index reads with illumina i5 and i7 index primers, respectively
106 (20+6+6 sequencing cycles).

107

108 **b. RSL guides allow additional methods of analysis.** In Total guide read Count Analysis (TCA, left) RSL
109 information is ignored and only the sum of readcounts for all RSL-guides is taken into account. In
110 internal replicate analysis (IRA, middle), readcounts of RSL-guides are binned such that internal
111 replicates are created for each guide. The example shown bins into four internal replicates. In
112 lineage dropout analysis (LDA, right) each RSL-guide is monitored separately.

113

114 **c. Screen size and sequencing depth.** The screens were performed at a very large screen size of
115 roughly 4500 cells per guide and sequenced to a depth of 30,000 reads per guide. Using RSL
116 information, the data from these oversized experiments were then subsampled bioinformatically

117 to approximately one quarter and one sixteenth, to test different analysis methods at different
118 screen sizes. The corresponding values for two published screens are indicated for comparison ^{4,5}
119
120 **d. Internal replicate analysis (IRA).** RSL-guides were binned to create 64 internal replicates. Effect
121 sizes (log₂ fold change in readcount between Day 4 and Day 28 after virus transduction) for each
122 bin are plotted in ascending order, 10 guides each for MYCN (**top left**) and MYC (**top right**), as well
123 as 50 representative non-targeting guides (**bottom**, these non-cutters seem to have a small
124 fitness advantage). Red dots, median effect size (MES) of the 64 internal replicates; black line, MES
125 of all guides in the library. Hits for this type of data were called from MES and SSMD
126 (**Supplementary Figure 1b**, see Supplementary Information for details).
127
128 **e. Lineage dropout analysis (LDA).** The average fraction of RSL-guides lost from day 4 to day 28 in
129 each experimental replicate is plotted for each gene. Red, positive controls; blue non-targeting
130 controls; black line, linear regression. The number of virus-transduced cell lineages lost is the most
131 direct readout of the guide effect on cell viability.
132

133 **Figure 2.**

134 **a. RSLs increase accuracy of hit calling.** Ranks of known positive controls (20 ribosomal proteins out
135 of a total of 2,335 interrogated genes) in one experimental replicate for the full screen size (left),
136 as well as one quarter (middle) and 1/16 (right) of the full screen size. Red line, median rank. At all
137 screen sizes, IR/SSMD analysis and LDA assigned lower ranks to the positive controls than TCA.
138 The variance of the ranking increased substantially with decreasing screen size in TCA, but not in
139 the two RSL-based methods.
140
141 **b. RSLs increase the precision of gene ranking.** Average percent overlap of the top-ranked 5% of
142 genes (116 genes) between two experimental replicates. Error bars, standard deviation of four
143 subsamples. LDA is the most precise method, followed by IRA/SSMD. Again, both RSL-based
144 methods are superior to TCA and much more robust at smaller screen sizes.
145
146 **c. RSLs boost statistical power at small screen size.** Hit gene overlap between experimental
147 replicates (<1% false discovery rate) at full screen size, one quarter and one sixteenth of the full
148 screen size. Error bars, standard deviation for hit gene overlap between four subsamples in
149 experimental replicate 1 and four subsamples in experimental replicate 2 (16 comparisons in
150 total). Only at full screen size, TCA matches LDA. At more practical screen sizes, LDA has much
151 higher statistical power and identifies considerably more hit genes.
152
153
154
155
156

157 **Accession codes**

158 Read data: European Nucleotide Archive, PRJEB18436. Scripts will be made available on Github under
159 public license.

160

161 **Acknowledgements**

162 The authors would like to thank Drs. Inderpreet Kaur Sur, Jenna Persson and Minna Taipale for suggestions
163 on the manuscript. Part of this work was carried out at Karolinska High Throughput Center (KHTC) and the
164 High Throughput Genome Engineering Facility (HTGE) funded by SciLifeLab.

165

166 **Authorship contributions.**

167 B.S., S.K.B. and J.T. developed the approach, B.S., S.K.B. and M.T. performed the experiments, B.S., S.K.B.,
168 J.Z. and T.K. analyzed the data, B.S. and J.T. wrote the manuscript.

169

170 **References**

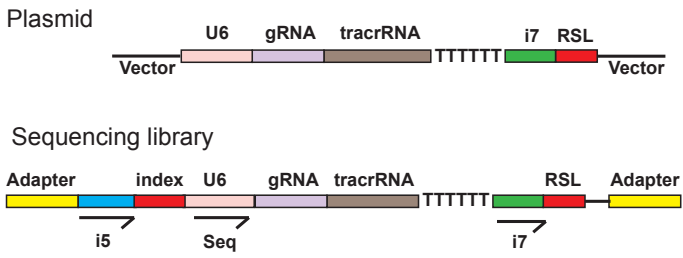
- 171 1. Shalem, O., Sanjana, N.E. & Zhang, F. High-throughput functional genomics using CRISPR-Cas9. *Nat*
172 *Rev Genet* **16**, 299-311 (2015).
- 173 2. Gilbert, L.A. et al. Genome-Scale CRISPR-Mediated Control of Gene Repression and Activation. *Cell*
174 **159**, 647-661 (2014).
- 175 3. Koike-Yusa, H., Li, Y., Tan, E.P., Velasco-Herrera Mdel, C. & Yusa, K. Genome-wide recessive genetic
176 screening in mammalian cells with a lentiviral CRISPR-guide RNA library. *Nat Biotechnol* **32**, 267-
177 273 (2014).
- 178 4. Shalem, O. et al. Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science* **343**, 84-
179 87 (2014).
- 180 5. Wang, T. et al. Identification and characterization of essential genes in the human genome. *Science*
181 **350**, 1096-1101 (2015).
- 182 6. Levy, S.F. et al. Quantitative evolutionary dynamics using high-resolution lineage tracking. *Nature*
183 **519**, 181-186 (2015).
- 184 7. Sandler, O. et al. Lineage correlations of single cell division time as a probe of cell-cycle dynamics.
185 *Nature* **519**, 468-471 (2015).
- 186 8. Kivioja, T. et al. Counting absolute numbers of molecules using unique molecular identifiers.
187 *Nature methods* **9**, 72-74 (2012).
- 188 9. Zhang, X.D. A pair of new statistical parameters for quality control in RNA interference high-
189 throughput screening assays. *Genomics* **89**, 552-561 (2007).
- 190 10. Li, W. et al. MAGeCK enables robust identification of essential genes from genome-scale
191 CRISPR/Cas9 knockout screens. *Genome Biol* **15**, 554 (2014).
- 192 11. Konermann, S. et al. Genome-scale transcriptional activation by an engineered CRISPR-Cas9
193 complex. *Nature* **517**, 583-588 (2015).

194

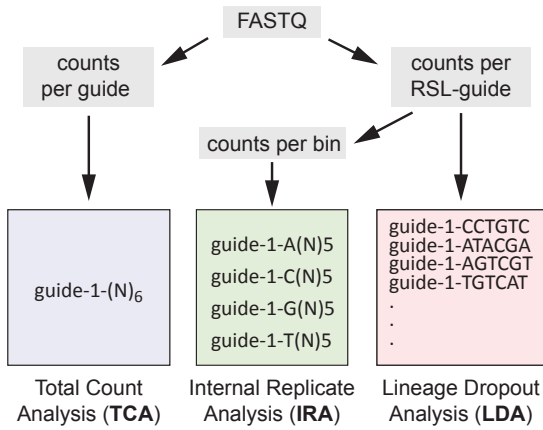
Figure 1

bioRxiv preprint doi: <https://doi.org/10.1101/114355>; this version posted May 15, 2017. 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.

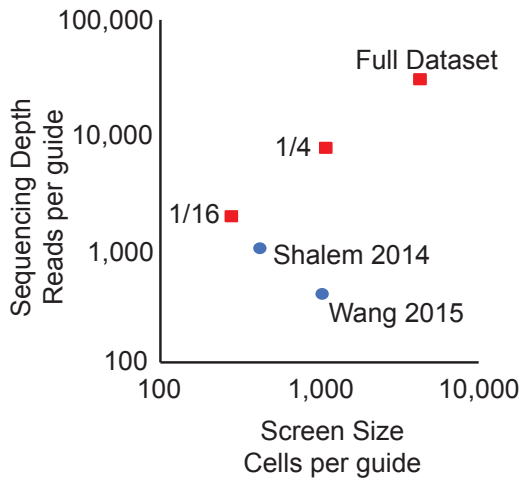
a. Library design



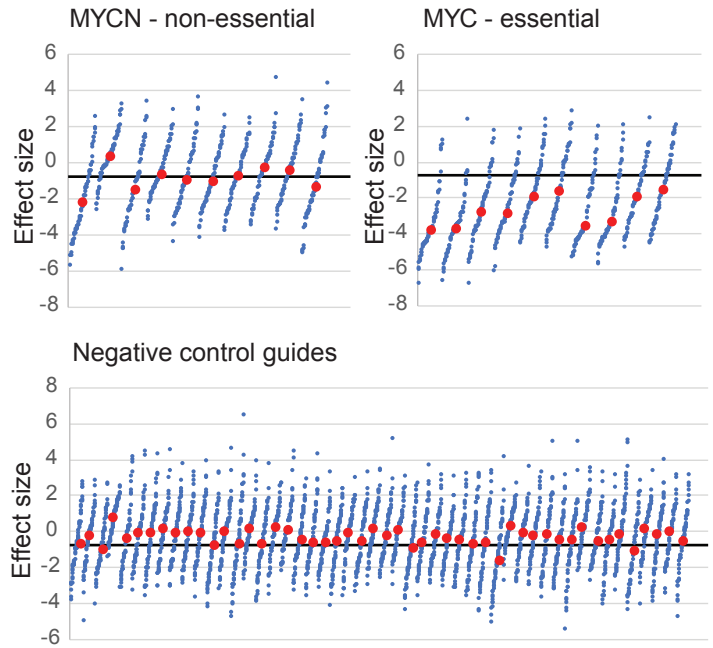
b. Levels of analysis



c. Screen size and sequencing depth



d. Internal replicate analysis (IRA)



e. Lineage dropout analysis (LDA)

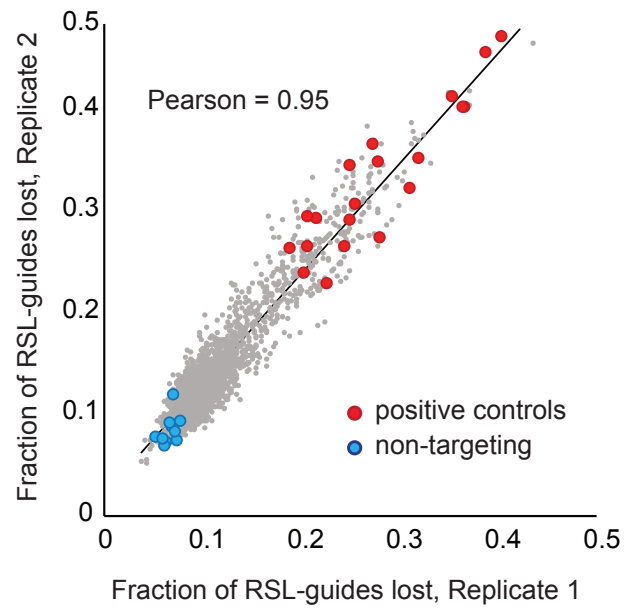
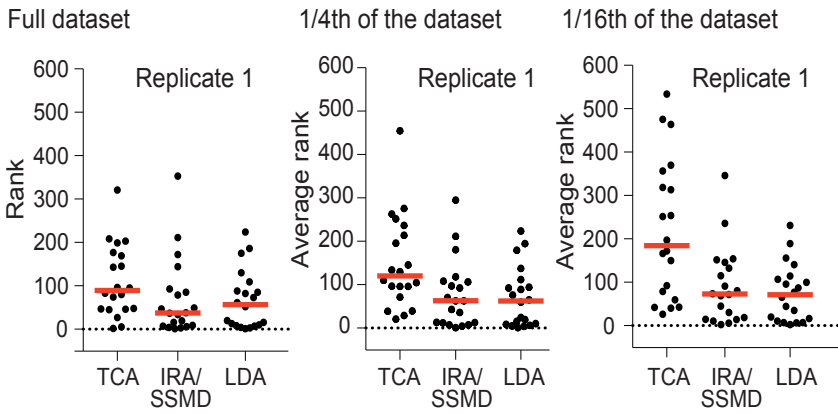


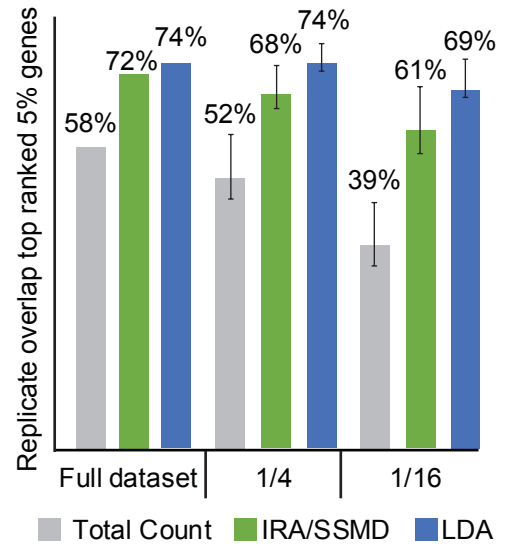
Figure 2

bioRxiv preprint doi: <https://doi.org/10.1101/114355>; this version posted May 15, 2017. 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.

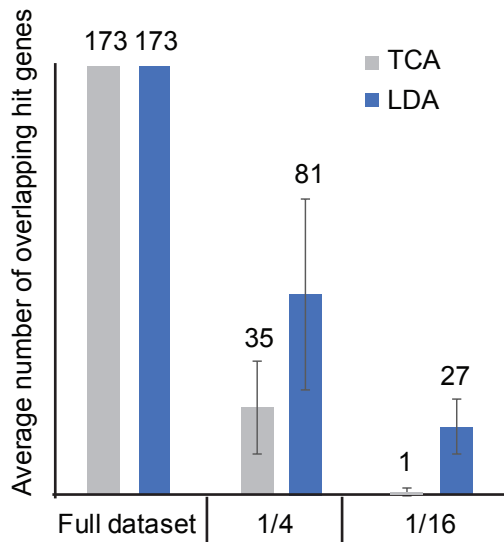
a. RSLs increase accuracy - Ranks of positive controls



b. RSLs increase precision - gene rank overlap



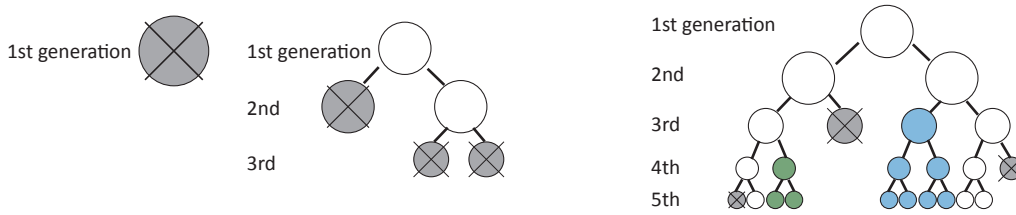
c. RSLs increase statistical power at smaller screen size



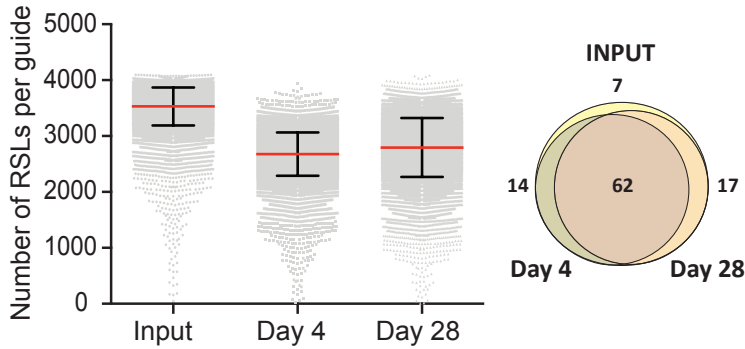
Schmierer et al., Supplementary Figure 1

a. Lineage drop out versus lineage depletion

bioRxiv preprint doi: <https://doi.org/10.1101/114355>; this version posted May 15, 2017. 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.

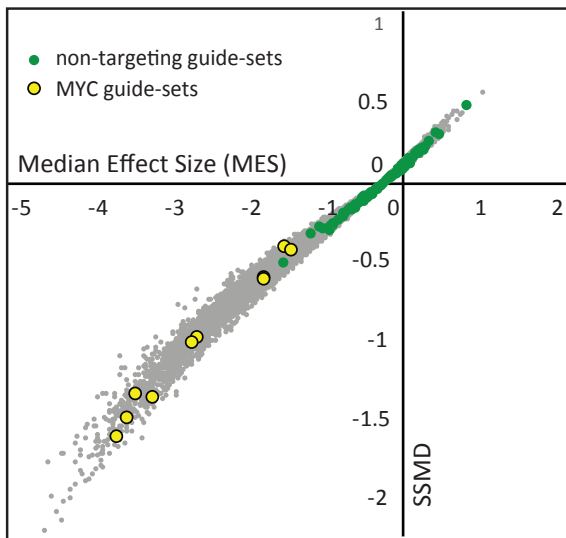


b. Library complexity and carry-through

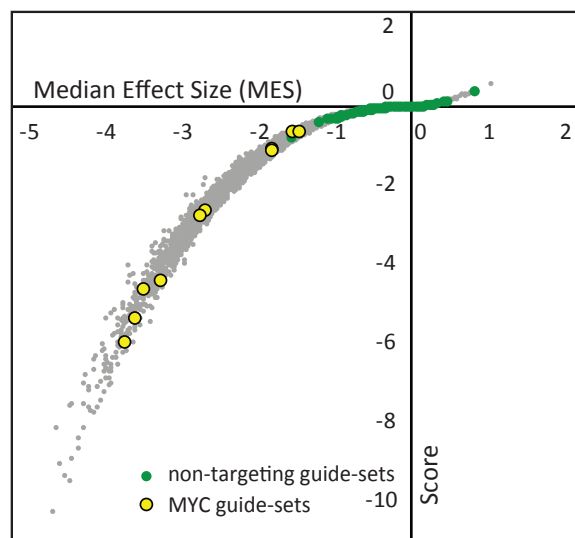


c. IRA analysis by median effect sized and SSMD

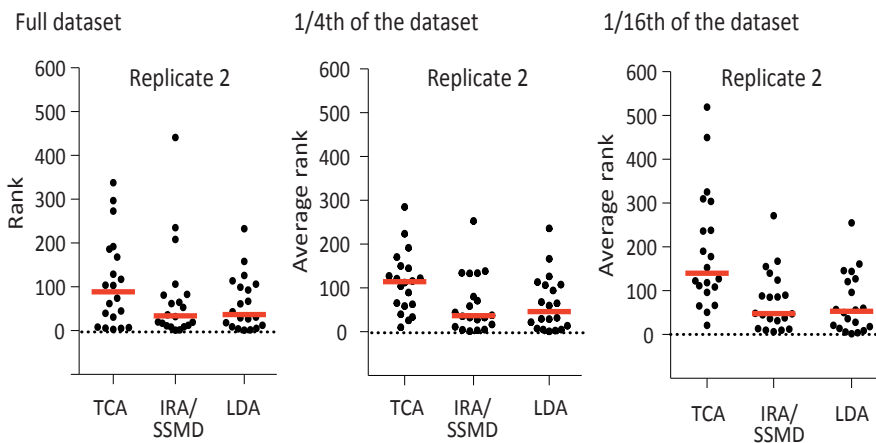
Double flashlight plot



IRA/SSMD Ranking Score



d. RSLs increase accuracy - Ranks of positive controls, replicate 2



Supplementary Figure S1.

- a. Lineage dropout versus lineage depletion.** Depending on the kinetics of editing, single cell lineages harboring a single RSL-guide against an essential gene can either disappear (drop out) or decrease in their abundance (depletion). **Left.** Dropout happens if the editing occurs early on, either before the cell can divide, or in several independent events at later time points (grey, dead cell; white, unedited cell). **Right.** In lineage depletion, editing occurs either after several cell divisions and/or with several different outcomes, some of which will retain gene function of the essential gene (blue and green edits). The traced lineage is then comprised of several sublineages.
- b. Number of distinct sequences carried through the experiment in one representative experimental replicate.** **Left.** Boxplot. An average of 3600 RSLs per guide in the plasmid library reduced to an average of 2800 RSLs per guide in the samples taken from the cell populations. Many of these RSLs have very low read counts, those were filtered out and not used for downstream data analysis. **Right.** Venn diagram. Both timepoints together covered 93% of input RSL-guides (78 million unique sequences in the cell population). The overlap between day 4 and day 28 was two thirds, with about one sixth of sequences found either only in Day 4 or only in Day 28. This is a consequence of unavoidable undersampling, which also occurs in the absence of RSLs, however in their presence becomes apparent and allows removal of inconsistently sampled lineages.
- c. Internal replicate analysis by strictly standardized mean difference (IRA/SSMD).** **Left.** Double flash-light plot. SSMDs for all >23,000 guide sets are plotted against their median effect sizes (MES). Green circles indicate non-targeting control guide sets, yellow circles indicate the 10 guide-sets targeting MYC. **Right.** Ranking score. A hit score was defined as the product of MES and SSMD. This score is negative for depleted guides and positive for enriched guides, and was used for guide ranking. The score was plotted against the median effect size for all guides. Green circles indicate non-targeting control guide sets, yellow circles indicate the 10 guide-sets targeting MYC.
- d. RSLs increase accuracy of hit calling.** As Fig. 2a, but for the second experimental replicate.

1 **Supplementary Information**

2

3 **Oligo synthesis and library cloning.** The guide library targeted 2325 genes and contains a total of
4 23,279 guides. The targeted gene set contains all human transcription factors¹, other genes of interest
5 as well as ribosomal proteins as positive controls and 101 non-targeting guides as negative controls.
6 All sgRNA sequences used in this library were taken from a previously published, genome-wide library²
7 (Supplementary file *RSL_guide_library.csv*). The 5' part of the library construct (blue + black, 122bp),
8 containing the sgRNA was synthesized by oligo array (CustomArray). The 3' part containing the RSL and
9 the Illumina i7 index primer sequence (green) was synthesized as a single 119 bp oligo
10 (black+green+red). These two oligos were annealed to each other at the overlapping part (black) and
11 double stranded by PCR ($T_M = 64C$) using outer primers (underlined).
12

13 GGCTTTATATATCTTGTGGAAAGGACGAAACACCGNNNNNNNNNNNNNNNNNNNNgtttAagagctag
14 aatagcaagttTaaataaggctagtcggttatcaacttgaaaaagtggcaccgagtcggtgcTTTTT
15 TgatcgggaagagcacacgtctgaactccagtcacNNNNNNaagcttggcgtaactagatcttgagaca
16 aa

17 The PCR product was cloned by Gibson assembly into the lentiviral vector pLenti-Puro-AU-flip-
18 3xBsMBI, which was created by modifying lentiGuide-Puro (a gift from Feng Zhang, Addgene #52963)
19 by replacing the sequence

20 gttttagagctagaaatagcaagttaaaataaggctagtcggttatcaacttgaaaaagtggcaccga
21 gtcggtgcTTTTTT

22 with

23 gtttAagagctagaaatagcaagttTaaataaggctagtcggttatcaacttgaaaaagtggcaccga
24 gtcggtgcTTTTTTCgtctct).

25

26 **Gibson assembly, transformation and amplification of the library.** 100 ng vector and 12 ng insert
27 where assembled in a total reaction volume of 100 μ l (NEBuilder[®] HiFi DNA Assembly Master Mix,
28 NEB). The reaction was cleaned via a Minelute reaction cleanup column (Qiagen) and transformed into
29 6 x 50 μ l electrocompetent *E. coli* (Endura[™] ElectroCompetent Cells, Lucigen) using a 1.0 mm cuvette,
30 25 μ F, 400 Ohms, 1800 Volts. Bacteria were plated on several 24x24 cm agar plates and colonies were
31 grown overnight. Colonies were scraped into LB medium and the contained plasmids were isolated by
32 Maxiprep.
33

34 **Library packaging.** The library was packaged in HEK 293T cells by cotransfecting the library plasmid
35 and the two packaging plasmids psPAX2 (a gift from Didier Trono, Addgene #12260) and pCMV-VSV-G
36 (a gift from Bob Weinberg, Addgene # 8454) in equimolar ratios. After 48 hours, the virus-containing
37 supernatant was concentrated 40-fold using Lenti-X concentrator (Clontech), aliquoted for one time
38 use and stored at -140C.
39

40 **Cell lines and cell culture.** All the cells used in this study were purchased directly from ATCC. Cells were
41 regularly tested for mycoplasma using the Mycoalert detection kit (Lonza; cat# LT07-218).
42

43 **Creating editing-proficient Cas9 cell lines.** To rapidly generate editing-proficient cell lines, we
44 synthesized a lentiviral construct (pLenti-Cas9-sgHPRT1) that encodes a codon optimized WT-SpCas9
45 that is flanked by two nuclear localization signals (derived from lenti-dCAS-VP64_Blast, a gift from Feng
46 Zhang, Addgene #61425). In addition, the construct codes for blasticidin resistance, and carries an
47 sgRNA against HPRT1 (GATGTGATGAAGGAGATGGG). HPRT1 loss confers resistance to the

48 antimetabolite 6-thioguanine (6-TG). Lentivirally transduced cells were selected in 5 µg/ml Blasticidin
49 and after one week to 10 days additionally with 5µg/ml 6-TG until control cells had died. Only cells that
50 both express Cas9 and are editing proficient, as indicated by loss of HPRT1 function, will survive. The
51 method allows rapid establishment of a pool of editing proficient cells. Compared to single cell clones,
52 this method retains the genetic heterogeneity of the original cell line, avoids potential clonal effects
53 of the particular integration site of Cas9, and greatly accelerates cell line generation. These benefits
54 need to be weighed carefully against possible disadvantages, such as synthetic lethality with HPRT1
55 loss, or potential effects of the presence of a second guide in the cell.

56
57 **Library transduction.** Per replicate, 100 million RKO Cas9 cells were transduced with the library virus.
58 Cells were then selected for guide integration and expression by 1 µg/ml puromycin selection for 48
59 hours. A proportion of cells will contain more than one guide. Because of the vast number of RSL-
60 guides, any ineffective passenger guides will associate with effective guides randomly and will not be
61 significantly enriched or depleted in the population.

62
63 **Cell propagation and sample preparation.** Cells were kept in culture for a total of 28 days after
64 transduction by sub-culturing them every three to four days. 100 million cells were reseeded at each
65 split, and genomic DNA was prepared from 50 – 80 million cells at Days 4 and 28 after transduction.
66 Day 4 after transduction was considered the control time point.

67
68 **Preparation of the sequencing library from genomic DNA.** The sequencing library preparation consists
69 of 3 PCR steps, PCR1 amplifies the genomic region containing the guide sequence using the primers 1F
70 and 1R. PCR2 and PCR3 then incorporate the Illumina adaptors with primers 2F/2R and 3F/3R,
71 respectively. 3F contains the Illumina index for multiplexing, indicated by NNNNNN in the sequence
72 given.

73 Genomic DNA was isolated using Blood and Tissue Maxi Kit (Qiagen), and 200 µg, theoretically
74 corresponding to 30 million diploid cells, were used as PCR template in 40 parallel PCR reactions (5 µg
75 template DNA each) using KAPA HiFi HotStart polymerase (KAPA Biosystems). After 14 cycles, the
76 reactions were pooled. PCR2 used 5 µl of pooled PCR1 as template and was run for 19 cycles, PCR3
77 used 2µl of PCR2 as template and was run for 14 cycles. The resulting product of 288 bp was gel purified
78 and sequenced on an Illumina HiSeq 4000 instrument using single read 20 cycles plus two 6 bp index
79 reads, where index read 1 reads the RSL and index read 2 reads the Illumina sample index.

80
81 1fw G GACTATCATATGCTTACCGTAACTTGAAAGTATTTTCG
82 1ref CTTTAGTTTGTATGTCTGTTGCTATTATGTCTACTATTCTTTCC
83 2fw TCTTCCCTACACGACGCTCTCCGATCtctgtggaagacgaaacac
84 2rev AGAAGACGGCATAACGAGATctgccattgtctcaagatctagttac
85 3fw AATGATACGGCGACCACCGAGATCTACAC NNNNNN TCTTCCCTACACGACGCTCTCCG
86 3rev CAAGCagaagacggcatacgagatctgccatttg

87
88 The final library product was sequenced with a custom primer and the i5 and i7 index primers
89 (underlined) by running 20+6+6 cycles on the Illumina HiSeq4000.

90
91 AATGATACGGCGACCACCGAGATCTACAC [i5] **NNNNNN**TCTTCCCTACACGACGCTCTCCGATCt
92 cttgtggaagacgaaacacCG**NNNNNNNNNNNNNNNNNNNN**gtttAagagctagaaatagcaagtt
93 TaaataaGgctagtcggttatcaacttgaaaaagtggcaccgagtcggtgcTTTTTTgatcggagag
94 cacacgtctgaactccagtcac [i7] **BBBBBB**aagcttggcgtaactagatcttgagacaaatggcag
95 ATCTCGTATGCCGCTTCTGCTTG

96
97

98 **Scripts used for counting RSL-guides and for binning.** RSL-guides were counted in the original fastq
99 files with the Perl scripts *BatchRun-pub2.pl*, which requires the script *GuideUMI_pub2.pl*. Binning of
100 RSL-guide counts was done using the script *countTruncatedRSLs.pl*. Generally, sequences whose total
101 readcount in control and treatment was less than five were filtered out prior to data analysis.

102

103 **SSMD analysis of read count data.**

104 **Normalization.** In RNASeq, methods such as median normalization are commonly preferred to total
105 read-count normalization, mainly to compensate for the effect of a few very highly expressed genes
106 that can take up a significant proportion of the total read count. CRISPR/Cas9 screening data are
107 comparably well balanced and we thus chose the most basic normalisation method, total read count
108 normalisation, to compensate for different sequencing depths. c_{ij} and t_{ij} represent the raw read
109 counts for RSL-guide j in guide-set i for control (Day 4 after lentiviral transduction) and treatment (Day
110 28 after lentiviral transduction), respectively. The normalised read counts c'_{ij} and t'_{ij} are then

111

$$112 \quad c'_{ij} = c_{ij} \frac{\sum_{ij}(c_{ij} + t_{ij})}{2 \sum_{ij} c_{ij}}$$

113

$$114 \quad t'_{ij} = t_{ij} \frac{\sum_{ij}(c_{ij} + t_{ij})}{2 \sum_{ij} t_{ij}}$$

115

116 **Median effect size and variability of the guide-sets.** We defined the effect size ES_{ij} for each RSL-guide
117 j in guide-set i as the log₂ of the fold change between treatment count and control count. To handle
118 total loss of an RSL-guide in the treatment sample, we added a pseudo-count of 1 to all counts:

119

$$120 \quad ES_{ij} = \log_2 \frac{t'_{ij} + 1}{c'_{ij} + 1}$$

121

122 Next, we calculated the median effect size for guide set i , MES_i , and the median of the absolute
123 deviations (MAD) of all RSL-guides j in guide-set i from MES_i

124

$$125 \quad MES_i = \text{median}_j ES_{ij}$$

126

$$127 \quad MAD_i = 1.4826 \text{median}_j |ES_{ij} - MES_i|$$

128

129 The factor 1.4826 was chosen such that the MAD is approximately equal to the standard deviation
130 under the assumption of normal distribution³.

131

132 **Median effect size and variability of the control guide-sets.** The RSL library contains 101 non-targeting
133 guide-sets. We calculate a single median effect size and MAD for this control set in the following way:

134

135 Median effect size of all non-targeting RSL-guides

136

$$137 \quad MES_{CON} = \text{median}_{ij} ES_{ij}^{NONT}$$

138

139 Median absolute deviation of all non-targeting RSL-guides:

140

$$141 \quad MAD_{CON} = 1.4826 \text{median}_{ij} |ES_{ij}^{NONT} - MES_{CON}|$$

142

143 **Strictly standardized mean difference (SSMD).** SSMD is a measure for the significance of the
144 difference in behaviour of sample i and the non-targeting controls. It takes into account both the effect
145 size and the variability of the data.

146

$$SSMD_i = \frac{MES_i - MES_{CON}}{\sqrt{MAD_i^2 + MAD_{CON}^2}}$$

147

148 For samples with relatively small effect size, the SSMD can still become large if the spread is small. We
149 thus introduce a score in which the effect size weighs more strongly, and which is used as a ranking
150 parameter:

151

$$Score_i = MES_i |SSMD_i|$$

152

153 For hit calling, the average score and standard deviation were calculated for all non-targeting guide
154 sets. The script used in these calculations is *SSMD.sh*, which calls the script R-script *SSMD.R*. Guide-
155 sets were then ranked according to their score and the resulting ranked list was analysed with α -RRA,
156 a robust rank aggregation algorithm as implemented in the “pathway” function of MAGeCK^{4,5} using
157 Supplementary file *RSL_guide_library.gmt*.

158

159

160

161

Lineage dropout.

162

163 An RSL-guide was considered a dropout if it had less than two readcounts in the treatment time point.
164 The numbers of RSLs per guide at Day 4 and Day 28 were then used as input in MAGeCK to obtain a
165 ranked gene list and FDRs.

166

167

Subsampling.

168

169 For subsampling the full data set, RSL-guides were grouped according to their RSL-sequence. For
170 medium screen size, the whole dataset was split into four groups (RSLs starting with A,C,G and T). For
171 small screen size, the whole dataset was split into 16 groups, the first four of which (AA,AC,AG,AT)
172 were used for analysis.

173

174

175

Supplementary References

176

- 177 1. Vaquerizas, J.M., Kummerfeld, S.K., Teichmann, S.A. & Luscombe, N.M. A census of human
178 transcription factors: function, expression and evolution. *Nat Rev Genet* **10**, 252-263 (2009).
- 179 2. Wang, T., Wei, J.J., Sabatini, D.M. & Lander, E.S. Genetic screens in human cells using the
180 CRISPR-Cas9 system. *Science* **343**, 80-84 (2014).
- 181 3. Zhang, X.D. Illustration of SSMD, z score, SSMD*, z* score, and t statistic for hit selection in
182 RNAi high-throughput screens. *J Biomol Screen* **16**, 775-785 (2011).
- 183 4. Kolde, R., Laur, S., Adler, P. & Vilo, J. Robust rank aggregation for gene list integration and
184 meta-analysis. *Bioinformatics* **28**, 573-580 (2012).
- 185 5. Li, W. et al. MAGeCK enables robust identification of essential genes from genome-scale
186 CRISPR/Cas9 knockout screens. *Genome Biol* **15**, 554 (2014).

187