

# 1 **CRISPulator: a discrete simulation tool for pooled genetic screens**

2

3 Tamas Nagy<sup>1</sup> and Martin Kampmann<sup>2\*</sup>

4

5 <sup>1</sup>Graduate program in Bioinformatics, University of California, San Francisco, CA 94158, USA

6 Tamas.Nagy@ucsf.edu

7

8 <sup>2</sup>Department of Biochemistry and Biophysics, Institute for Neurodegenerative Diseases

9 California Institute for Quantitative Biomedical Research, University of California, and Chan-

10 Zuckerberg Biohub, San Francisco CA 94158, USA

11 Martin.Kampmann@ucsf.edu

12 \*Corresponding Author

13

## 14 **Abstract**

15 The rapid adoption of CRISPR technology has enabled biomedical researchers to conduct  
16 CRISPR-based genetic screens in a pooled format. The quality of results from such screens is  
17 heavily dependent on optimal screen design, which also affects cost and scalability. We present  
18 CRISPulator, a computational tool that simulates the impact of screen parameters on the  
19 robustness of screen results, thereby enabling users to build intuition and insights that will  
20 inform their experimental strategy. We illustrate its power by deriving non-obvious rules for  
21 optimal screen design.

22

## 23 **Keywords**

24 CRISPR; CRISPRi; functional genomics; genome-wide screens; simulation;

25

## 26 **Background**

27 Genetic screening is a powerful discovery tool in biology that provides an important functional  
28 complement to observational genomics. Until recently, screens in mammalian cells were  
29 implemented primarily based on RNA interference (RNAi) technology. Inherent off-target  
30 effects of RNAi screens present a major challenge [1]. In principle, this problem can be  
31 overcome using optimized ultra-complex RNAi libraries [2, 3], but the resulting scale of the  
32 experiment in terms of the number of cells required to be screened can be prohibitive for some  
33 applications, such as screens in primary cells or mouse xenografts.

34         Recently, several platforms for mammalian cell screens have been implemented based on  
35 CRISPR technology [4]. CRISPR nuclease (CRISPRn) screens [5, 6] perturb gene function by  
36 targeting Cas9 nuclease programmed by a single guide RNA (sgRNA) to a genomic site inside  
37 the coding region of a gene of interest, followed by error-prone repair through the cellular non-  
38 homologous end-joining pathway. CRISPR interference (CRISPRi) and CRISPR activation  
39 (CRISPRa) screens [7] repress or activate the transcription of genes by exploiting a catalytically  
40 dead Cas9 to recruit transcriptional repressors or activators to their transcription start sites, as  
41 directed by sgRNAs.

42         CRISPRn and CRISPRi have vastly reduced off-target effects compared with RNAi, and  
43 thus overcome a major challenge of RNAi-based screens. However, other challenges to  
44 successful screening [1] remain. The majority of CRISPRi and CRISPRn screens have been  
45 carried out as pooled screens with lentiviral sgRNA libraries. While this pooled approach has  
46 enabled rapid generation and screening of complex libraries, successful implementation of

47 pooled screens requires careful choices of experimental parameters. Choices for many of these  
48 parameters represent a trade-off between optimal results and cost.

49

## 50 **Results**

51 Here, we present a computational tool, termed CRISPulator, which simulates how experimental  
52 parameters will affect the detection of different types of gene phenotypes in pooled CRISPR-  
53 based screens. CRISPulator is freely available online (<http://crispulator.ucsf.edu>) to enable  
54 researchers to develop an intuition for the impact of experimental parameters on pooled  
55 screening results, and to optimize the design of pooled screens for specific applications. It  
56 simulates all steps of pooled screens, as visualized in **Fig. 1** and described in more detail in the  
57 Methods.

58 Briefly, a theoretical genome is generated in which genes are assigned quantitative  
59 phenotypes (**Fig. 2**). Independently, the quantitative relationship between gene knockdown level  
60 and resulting phenotype is defined for each gene (**Fig. 3**). Next, a sgRNA library targeting this  
61 genome is defined. Each gene is targeted by a number of independent sgRNAs. The technical  
62 performance of each sgRNA is randomly assigned based on a user-defined distribution of  
63 CRISPRn or CRISPRi sgRNA activities (**Fig. 4**), and the initial frequency distribution is  
64 specified (**Fig. 5**).

65 Simulation of the screen itself discretely models infection of cells with the pooled sgRNA  
66 library, phenotypic selection of cells and quantification of sgRNA frequencies in selected cell  
67 populations by next-generation sequencing. Based on the resulting data, hit genes are called (**Fig.**  
68 **6**) using our previously described quantitative framework [3], as detailed in the Online Methods.  
69 The performance of the screen with a specific set of experimental parameters is evaluated by

70 comparing the called hit genes to the actual genes with phenotypes defined by the theoretical  
71 genome. It is quantified either as overlap of the list of top called hits with the actual list of top  
72 hits, or as area under the precision-recall curve (AUPRC), a metric commonly used in machine  
73 learning (**Fig. 7**).

74 A central consideration for all pooled screens is the number of cells used relative to the  
75 number of different sgRNAs in the library. We refer to this parameter as representation, and  
76 distinguish representation at the time of infection, representation at times during phenotypic  
77 selection, and – by extension – representation at the sequencing stage (where it is defined as the  
78 number of sequencing reads relative to the relative to the number of different sgRNAs). From  
79 first principles, higher representation is desirable to reduce Poisson sampling noise (“jackpot  
80 effects”); in practical terms, higher representation is also more costly. A major application of  
81 CRISPulator is the exploration of parameters to guide the choice of suitable representation at  
82 each step of the screen to enable researchers to strike the desired balance between screening cost  
83 and performance.

84 CRISPulator implements two distinct strategies for phenotypic selection. In fluorescence-  
85 activated cells sorting (FACS)-based screens, cell populations are separated based on a  
86 fluorescent reporter signal that is a function of the phenotype. We [8] and others [9] have  
87 successfully implemented such screens by isolating and comparing cell populations with the  
88 highest and the lowest reporter levels. More commonly, pooled screens are conducted to detect  
89 genes with growth or survival phenotypes [5-7] by comparing cell populations at an early time  
90 point with cells grown in the absence or presence of selective pressures, such as drugs or toxins.

91 We first asked how representation at the infection, selection and sequencing stages  
92 affects FACS- and growth-based screens (**Fig. 8**). The performance of FACS-based screens was

93 most sensitive to the representation at the selection bottleneck, and least sensitive to  
94 representation at the infection stage, highlighting the importance of collecting a sufficient  
95 number of cells for each population during FACS sorting, ideally more than 100-fold the number  
96 of different library elements. By contrast, the performance of growth-based screens was similarly  
97 sensitive to representation at all stages.

98 For FACS screens using a given number of cells, an important decision is how extreme  
99 the cutoffs defining the “high-reporter” and “low-reporter” bins should be. CRISPulator  
100 simulation suggests that separating and comparing the cells with the top quartile and bottom  
101 quartile reporter activity results in the optimal detection of hit genes (**Fig. 9**). Closer inspection  
102 revealed that while both signal (sgRNA frequency differences between the two populations) and  
103 the noise (due to lower representation in the sorted population) decrease with larger bin sizes, the  
104 signal-to-noise ratio reaches a local maximum around 25% (**Fig. 10**), close to the bin size chosen  
105 fortuitously in published studies [8, 9].

106 For growth-based screens, the duration of the screen influences the signal (by amplifying  
107 differences in frequency due to different growth phenotypes) but also the noise (by increasing the  
108 number of Poisson sampling bottlenecks generated by cell passaging or repeated applications of  
109 selective pressure). Interestingly, CRISPulator suggests that the effect of screen duration on  
110 optimal performance is different for genes with positive and negative phenotypes, and strongly  
111 depends on the presence of genes with positive phenotypes (**Fig. 11**). While genes with positive  
112 phenotypes (increased growth / survival) were detected more reliably after longer screens, genes  
113 with negative phenotypes (decreased growth / survival) were optimally detected in screens of  
114 intermediate duration, and their detection in longer screens rapidly declined if genes with  
115 stronger positive phenotypes were present in the simulated genome. While genes with positive

116 phenotypes are rare in screens based on growth in standard conditions [5-7], selective pressures,  
117 such as growth in the presence of toxin, can reveal strong positive phenotypes for genes  
118 conferring resistance to the selective pressure [7]. The optimal screen length for growth-based  
119 screens was dictated by a local maximum of the signal-to-noise ratio, which itself depended on  
120 the representation: screens with lower representation were performing better at shorter duration  
121 (**Fig. 12**). Our results therefore predict that especially for growth-based screens using selective  
122 pressures, and screens implemented with low representation, short durations are preferable.

123 While CRISPRn and CRISPRi screens performed similarly in the simulations described  
124 above (**Fig. 8-11**), separate evaluation of genes with linear versus sigmoidal phenotype-  
125 knockdown relationship revealed that CRISPRn outperforms CRISPRi for the detection of  
126 sigmoidal genes (which require very stringent knockdown to result in a phenotype), whereas  
127 CRISPRi performs relatively better for genes with a linear knockdown-phenotype relationship  
128 (**Fig. 13**).

## 129

## 130 **Discussion**

131 CRISPulator revealed several non-obvious rules for the design of pooled genetic screens,  
132 illustrating its usefulness. Since certain parameters used by CRISPulator (such as the quality of  
133 sgRNA libraries or the signal-to-noise of FACS-based phenotypes) are estimates informed by  
134 published data, but not directly known, the predicted screen performance does not represent  
135 absolute performance metrics. Rather, the goal is to predict the relative performance of screens  
136 conducted with different experimental parameters to enable researchers to optimize those  
137 parameters. The simulated sequencing reads generated by CRISPulator (**Fig. 10**) recapitulate  
138 patterns observed in experimental data (**Fig. 14**), thereby facilitating the interpretation of

139 suboptimal experimental data and providing a tool to predict which experimental parameters  
140 need to be changed to obtain data more suitable for robust hit detection

141

## 142 **Conclusions**

143 CRISPulator facilitates the design of pooled genetic screens by enabling the exploration of a  
144 large space of experimental parameters *in silico*, rather than through costly experimental trial and  
145 error. For pooled genetic screens in animal models, such as mice, choices of experimental  
146 parameters can also have ethical implications, namely the numbers of animals required to power  
147 the study. As larger numbers of pooled genetic screens are published, we will further refine the  
148 assumptions underlying the simulation using empirical data.

149

## 150 **Methods**

151

152 **Code implementation and availability.** CRISPulator was implemented in Julia  
153 (<http://julialang.org>), a high-level, high-performance language for technical computing. We have  
154 released the simulation code as a Julia package, Crispulator.jl. The software is platform-  
155 independent and is tested on Linux, OS X (macOS), and Windows. Installation details,  
156 documentation, source code, and examples are all publicly available at  
157 <http://crispulator.ucsf.edu>.

158

159 **Simulated genome.** A genome is defined by assigning a numerical, “true” phenotype to a  
160 number of genes. All of our results featured here used 500 genes in each simulation. 75% of  
161 genes were assigned a phenotype of 0 (wild-type), and 5% of genes were modeled as negative

162 control genes, also with a phenotype of 0. 10% of genes were assigned a positive phenotype  
163 randomly drawn (unless otherwise indicated) from a Gaussian distribution with  $\mu=0.55$  and  
164  $\sigma=0.2$  (clamped between  $[0.1, 1.0]$ ), and 10% of genes were assigned a negative phenotype  
165 randomly drawn from an identical distribution except with  $\mu=-0.55$  and clamping  $[-1.0, -0.1]$   
166 (**Fig. 2**). Next, each gene was randomly assigned a phenotype-knockdown function (**Fig. 3**) to  
167 simulate different responses of genes to varying levels of knockdown. 75% of genes were  
168 assigned a linear function that linearly interpolates between 0 and the “true” phenotype from  
169 above as a function of knockdown, the remaining 25% of genes were assigned a sigmoidal  
170 function with an inflection point,  $p$ , drawn from a distribution with a mean of 0.8 and standard  
171 deviation of 0.2; the width of the inflection region,  $k$ , (over which a phenotype increased from 0  
172 to the “true” phenotype,  $l$ ) was drawn from a normal distribution with a mean of 0.1 and a  
173 standard deviation of 0.05. The function  $f$  was defined as follows:

174

$$f(x) = \begin{cases} 0, & x \leq p - k \\ 1, & x \geq p + k \\ \frac{1}{2} \left( \frac{\text{sign}(\delta) \cdot 1.05 |\delta|}{|\delta| + 1} + 1 \right), & p - k < x < p + k \end{cases}$$

175

176 where  $\delta = \frac{x-p}{\min(p, \min(1-p, k))}$

177

178 This specific sigmoidal function was chosen over the more standard special case of the logistic  
179 function or the Gompertz function because it is highly tunable and has a range between 0 and  $l$   
180 on a domain of  $[0, 1]$ .

181



182 **Simulated sgRNA libraries.** CRISPRn and CRISPRi sgRNA libraries are generated to target  
183 the simulated genome. For the results featured here, each gene was targeted by 5 independent  
184 sgRNAs. For CRISPRi screens, each sgRNA was randomly assigned a knockdown efficiency  
185 from a bimodal distribution (**Fig. 4**): 10% of sgRNAs had low activity with a knockdown drawn  
186 from a Gaussian ( $\mu=0.05$ ,  $\sigma=0.07$ ), 90% of guides had high activity drawn from a Gaussian  
187 ( $\mu=0.90$ ,  $\sigma=0.1$ ). We assumed such a high rate of active sgRNAs based on our recently  
188 developed highly active CRISPRi sgRNA libraries [10]. For CRISPRn screens, high-quality  
189 guides all had a maximal knockdown efficiency of 1.0 and were 90% of the population (the 10%  
190 low-activity CRISPRn guides were drawn from the same Gaussian ( $\mu=0.05$ ,  $\sigma=0.07$ ) as above).  
191 The initial frequency distribution of sgRNAs in the library was modeled as a log-normal  
192 distribution such that a guide in the 95<sup>th</sup> percentile of frequencies is 10 times as frequent as one  
193 in the 5<sup>th</sup> percentile (**Fig. 5**), which is typical of high-quality libraries in our hands [7].

194  
195 **Simulated screens.** All steps of the pooled screens are simulated discretely. Infections are  
196 modeled as a Poisson process with  $\lambda=M.O.I$  of the infection. The initial pool of cells is randomly  
197 infected by sgRNAs based on the frequency of each sgRNA in the library. A  $\lambda=0.25$  is used  
198 unless otherwise noted, which is commonly used to approximate single-copy infection [11]. Only  
199 cells with a single sgRNA are then used in subsequent steps, which is  $P(x=1; \text{Poisson}(\lambda=0.25)) \approx$   
200 19.5% of the initial pool.

201 For CRISPRi screens, phenotypes for each cell were determined based on the sgRNA  
202 knockdown efficiency (from above) and based on both the phenotype and the knockdown-  
203 phenotype relationship of the targeted gene. For CRISPRn screens, phenotypes for each cell  
204 were set using using sgRNA knockdown efficiency (specific for CRISPRn screens, see previous

205 paragraph) and the gene phenotype. If a cell was infected with a low-quality CRISPRn guide, it  
206 behaved similarly to a low-quality CRISPRi guide, i.e. very close to no activity. All cells with  
207 high-quality guides CRISPRn guides had a 1/9, 4/9, or 4/9 chance of having 0%, 50%, or 100%  
208 knockdown efficiency, respectively. This knockdown efficiency was then used with the  
209 knockdown-phenotype relationship and true phenotype of the gene to calculate the observed  
210 phenotype. The assumption that only bi-allelic frame-shift mutations lead to a phenotype in  
211 CRISPRn screens for most sgRNAs is supported by the empirical finding that in-frame deletions  
212 mostly do not show strong phenotypes, unless they occur in regions encoding conserved residues  
213 or domains [10]. To mitigate this issue, some CRISPRn screens have been conducted in quasi-  
214 haploid cell lines [6].

215 FACS sorting was simulated by convoluting the theoretical phenotypes of each cell  
216 independently with a Gaussian ( $\mu=0$ ,  $\sigma$ ) where  $\sigma$  is a tunable “noise” parameter, reflecting  
217 biological variance in fluorescence intensity of isogenic cells. The number of cells prior to this  
218 step is termed the bottleneck representation and is tunable. Post-convolution, cells were sorted  
219 according to their new, “observed” phenotype and then the bottom X percentile and top X  
220 percentile (X was real value between 0 and 50) were taken as the two comparison bins.

221 Growth experiments were simulated as follows: (1) in the time frame that WT cells (true  
222 phenotype=0) divide once, cells with the maximal negative phenotype, -1, do not divide, and  
223 cells with maximal positive phenotype divide twice. For cells with phenotypes in between 0 and  
224  $\pm 1$ , cells randomly pick whether they behave like WT cells or maximal phenotype cells weighted  
225 by their phenotype (i.e. cells with phenotypes close to 0 behave mostly like WT cells). (2) After  
226 one timestep where WT cells double once, a random subsample of the cells is taken. The size of

227 the bottleneck is tunable. (3) This is repeated  $n$  number of times. Finally, the samples of cells at  
228  $t=0$  and  $t=n$  are taken as the two populations for comparison.

229  
230 Sample preparation was simulated by taking the frequencies of each guide in the cells  
231 after selection and constructing a categorical distribution with the frequencies as the weights.  
232 Next-generation sequencing was then simulated by sampling from this categorical distribution up  
233 to the number of total reads.

234  
235 **Evaluation of screen performance.** Based on the simulated sequencing read counts, P values  
236 and gene-level phenotypes were calculated for each gene essentially as previously described [3,  
237 7]. Briefly, observed sgRNA phenotypes were calculated as  $\log_2$  ratios of sgRNA frequencies in  
238 two cell populations. Gene-level phenotypes were calculated by averaging the sgRNA  
239 phenotypes. P values were calculated based on the Mann-Whitney rank-sum test by comparing  
240 the phenotypes of sgRNAs targeting a given gene with the phenotypes of negative control  
241 sgRNAs. Genes were ranked by the product of the absolute gene-level phenotype and their –  
242  $\log_{10}$  P value to call hit genes. Screen performance was quantified in two ways: As the overlap of  
243 the top 50 called hit genes with the top 50 actual hit genes (based on true phenotype), or as the  
244 area under the precision-recall curve (AUPRC). AUPRC was chosen over the more common area  
245 under the receiver operator characteristic (AUROC) due to the highly skewed nature of the  
246 generated dataset (<20% of dataset is made up of true hits). The AUPRC was calculated using a  
247 lower trapezoidal estimator, which had been previously shown to be a robust estimator of the  
248 metric [12]. For **Figures 10** and **12**, the “signal” of an experiment was defined as the median  
249 signal for true hit genes (ones initially labeled as having a positive or negative phenotype). The

250 true hit gene signal was calculated as the average ratio of the  $\log_2$  fold change over the  
251 theoretical phenotype of all guides targeting that gene. Guides that dropped out of the analysis  
252 were excluded from the signal calculation. “Noise” was quantified as the standard deviation of  
253 negative-control sgRNA phenotypes, and the “signal-to-noise” ratio was the ratio of these two  
254 metrics. For display purposes, all are normalized in each graph.

255

### 256 **List of abbreviations**

257 AUPRC, area under the precision-recall curve

258 CRISPRi, CRISPR interference

259 CRISPRn, CRISPR nuclease

260 FACS, fluorescence-activated cell sorting

261 sgRNA, single guide RNA

262

### 263 **Declarations**

264

### 265 **Ethics approval and consent to participate**

266 Not applicable.

267

### 268 **Consent for publication**

269 Not applicable.

270

271 **Availability of data and materials**

272 The software, CRISPulator, described in this study, which was also used to generate the data, is  
273 publicly available at <http://crispulator.ucsf.edu> (an archived version is available at  
274 <https://doi.org/10.5281/zenodo.345188>). It is released under Apache License 2.0 and there are no  
275 additional restrictions for commercial use.

276

277 **Competing interest**

278 MK is an inventor on a patent application related to CRISPRi and CRISPRa screening  
279 (PCT/US15/40449).

280

281 **Funding**

282 TN was supported by an NSF graduate research fellowship. MK was supported by NIH/NIGMS  
283 New Innovator Award DP2 GM119139. The funding bodies had no role in the design of the  
284 study and collection, analysis, and interpretation of data and in writing the manuscript.

285

286 **Authors' contributions**

287 TN and MK designed and analyzed the research and wrote the manuscript. TN developed the  
288 CRISPulator software.

289

290 **Acknowledgments**

291 We thank Ruilin Tian and Daniel Asarnow for feedback on this manuscript, and Diane  
292 Nathaniel, John Chen, Kathleen Keough and Xiaoyan Guo for sharing unpublished experimental  
293 data.

294

295 **References**

296

- 297 1. Kaelin WG, Jr.: **Molecular biology. Use and abuse of RNAi to study mammalian gene**  
298 **function.** *Science* 2012, **337**:421-422.
- 299 2. Kampmann M, Horlbeck MA, Chen Y, Tsai JC, Bassik MC, Gilbert LA, Villalta JE,  
300 Kwon SC, Chang H, Kim VN, Weissman JS: **Next-generation libraries for robust**  
301 **RNA interference-based genome-wide screens.** *Proc Natl Acad Sci U S A* 2015,  
302 **112**:E3384-E3391.
- 303 3. Kampmann M, Bassik MC, Weissman JS: **Integrated platform for genome-wide**  
304 **screening and construction of high-density genetic interaction maps in mammalian**  
305 **cells.** *Proc Natl Acad Sci U S A* 2013, **110**:E2317-2326.
- 306 4. Shalem O, Sanjana NE, Zhang F: **High-throughput functional genomics using**  
307 **CRISPR-Cas9.** *Nat Rev Genet* 2015, **16**:299-311.
- 308 5. Shalem O, Sanjana NE, Hartenian E, Shi X, Scott DA, Mikkelsen TS, Heckl D, Ebert  
309 BL, Root DE, Doench JG, Zhang F: **Genome-scale CRISPR-Cas9 knockout screening**  
310 **in human cells.** *Science* 2014, **343**:84-87.
- 311 6. Wang T, Wei JJ, Sabatini DM, Lander ES: **Genetic screens in human cells using the**  
312 **CRISPR-Cas9 system.** *Science* 2014, **343**:80-84.
- 313 7. Gilbert LA, Horlbeck MA, Adamson B, Villalta JE, Chen Y, Whitehead EH, Guimaraes  
314 C, Panning B, Ploegh HL, Bassik MC, et al: **Genome-Scale CRISPR-Mediated Control**  
315 **of Gene Repression and Activation.** *Cell* 2014, **159**:647-661.

- 316 8. Sidrauski C, Tsai JC, Kampmann M, Hearn BR, Vedantham P, Jaishankar P, Sokabe M,  
317 Mendez AS, Newton BW, Tang EL, et al: **Pharmacological dimerization and**  
318 **activation of the exchange factor eIF2B antagonizes the integrated stress response.**  
319 *Elife* 2015, **4**:e07314.
- 320 9. DeJesus R, Moretti F, McAllister G, Wang Z, Bergman P, Liu S, Frias E, Alford J,  
321 Reece-Hoyes JS, Lindeman A, et al: **Functional CRISPR screening identifies the**  
322 **ufmylation pathway as a regulator of SQSTM1/p62.** *Elife* 2016, **5**.
- 323 10. Horlbeck MA, Gilbert LA, Villalta JE, Adamson B, Pak RA, Chen Y, Fields AP, Park  
324 CY, Corn JE, Kampmann M, Weissman JS: **Compact and highly active next-**  
325 **generation libraries for CRISPR-mediated gene repression and activation.** *Elife*  
326 2016, **5**.
- 327 11. Fellmann C, Zuber J, McJunkin K, Chang K, Malone CD, Dickins RA, Xu Q, Hengartner  
328 MO, Elledge SJ, Hannon GJ, Lowe SW: **Functional identification of optimized RNAi**  
329 **triggers using a massively parallel sensor assay.** *Mol Cell* 2011, **41**:733-746.
- 330 12. Boyd K, Eng KH, Page CD: In *Machine Learning and Knowledge Discovery in*  
331 *Databases*. Edited by Blockeel H, Kersting K, Nijssen S: Springer; 2013: 451-466

332

333

### 334 **Figure legends**

335 **Figure 1.** CRISPulator simulates pooled genetic screens to evaluate the effect of experimental  
336 parameters on screen performance. Overview of simulation steps: Parameters listed with bullet  
337 points can be varied to examine consequences on the performance of the screen, which is

338 evaluated as the detection of genes with phenotypes (quantified as overlap or area under the  
339 precision-recall curve, AUPRC). Details are given in the text and Methods.

340 **Figure 2.** Phenotype distribution in the simulated genome. A typical distribution is shown, which  
341 includes 75% of genes without phenotype (green), 5% of negative control genes (pink), 10% of  
342 genes with a positive phenotype (blue), and 10% of genes with a negative phenotype (yellow).

343

344 **Figure 3.** Relationship between gene knockdown level and resulting phenotype. This  
345 relationship is defined for each gene, and represents either a linear function (orange graphs) or a  
346 sigmoidal function (blue lines), as defined in the Online Methods.

347

348 **Figure 4.** CRISPRi sgRNA activity distribution. An example of a typical distribution for 1000  
349 guides is shown.

350

351 **Figure 5.** Initial frequency distribution of sgRNAs. An example of a typical distribution is  
352 shown.

353

354 **Figure 6.** Sample results from a CRISPulator simulation of a FACS-based screen. Top row:  
355 Each point represents an individual sgRNA, plotting its read numbers in the simulated deep  
356 sequencing run for the “low reporter signal” bin and the “high reporter signal” bin. sgRNAs are  
357 color-coded to indicate whether they target a gene with a positive phenotype (knockdown  
358 increases reporter signal, blue), a gene with a negative phenotype (knockdown decreases reporter  
359 signal, red), a gene without phenotype (grey), or whether they are non-targeting control sgRNAs  
360 (black). Bottom row: Based on the observed sgRNA phenotypes, gene phenotypes are calculated



361 (mean  $\log_2$  ratio of read frequencies in “high” over “low” bins), and a gene P value is calculated  
362 to express statistical significance of deviation from wild-type. These are visualized in volcano  
363 plots in which each dot represents a gene. Genes are color-coded to indicate the actual  
364 phenotype: positive, blue; negative, red; no phenotype, grey.

365  
366 **Figure 7.** Metrics to evaluate screen performance. (a) “Venn diagram” overlap between the 50  
367 genes with the strongest actual phenotypes, and the top 50 hit genes called based on the screen  
368 results – expressed as the ratio of the number of genes in the overlap over the number of called  
369 top hit genes, i.e. 50. (b) Area under the precision-recall curve (AUPRC).

370  
371 **Figure 8.** Importance of representation of library elements at different stages of the screen.  
372 CRISPulator simulations reveal the effect of library representation at different screen stages  
373 (Transfection, bottlenecks, sequencing) on hit detection. Simulations were run for FACS-based  
374 screens (top row) and growth-based screens (bottom row). Lines and light margins represent  
375 means and 95% confidence intervals, respectively, for 10 independent simulation runs.

376  
377 **Figure 9.** Effect of bin size on performance of FACS-based screens. Simulations were run for  
378 100x representation at the transfection, bottleneck and sequencing stages. Lines and light  
379 margins represent means and 99% confidence intervals, respectively, for 100 independent  
380 simulation runs.

381  
382 **Figure 10.** Effect of bin size on signal and noise of FACS-based screens. For FACS-based  
383 screens, the effect of the size of the sorted bins (see Fig. 1) on metrics for signal, noise, and

384 signal-to-noise ratio (scaled within each plot) is shown. Metrics are defined in the Online  
385 Methods. Simulations were run for 10x representation (top row) or 100x representation (bottom  
386 row) at the transfection, bottleneck and sequencing stages. Lines and light margins represent  
387 means and 99% confidence intervals, respectively, for 25 independent simulation runs.

388  
389 **Figure 11.** Effect of positive phenotypes on growth-based screens. For growth-based screens, the  
390 presence of genes with positive phenotypes (fitter than wild type) strongly influences hit  
391 detection as a function of screen duration. Screens were simulated for a set of genes in which  
392 10% of all genes had negative phenotypes (less fit than wild type), and 2% of genes had positive  
393 phenotypes. The strength of positive phenotypes was varied, as encoded by the heat map. Hit  
394 detection was quantified separately for genes with negative phenotypes (top row) and genes with  
395 positive phenotypes (bottom row). Simulations were carried out for screens with different  
396 durations, as measured by the number of passages. Lines and light margins represent means and  
397 95% confidence intervals, respectively, for 25 independent simulation runs. In **a** and **c**, hit  
398 detection is measured as Area under the Precision-Recall curve (AUPRC), as detailed in the  
399 Online Methods.

400  
401 **Figure 12.** Effect of duration of growth-based screens on performance. Screens were simulated  
402 for a set of genes in which 10% of all genes had negative phenotypes (less fit than wild type).  
403 Simulations were carried out for screens with different durations, as measured by the number of  
404 passages, and for different representations at the transfection, bottleneck and sequencing stages.  
405 Metrics for signal, noise, and signal-to-noise ratio are defined in the Online Methods. Lines and

406 light margins represent means and 95% confidence intervals, respectively, for 25 independent  
407 simulation runs.

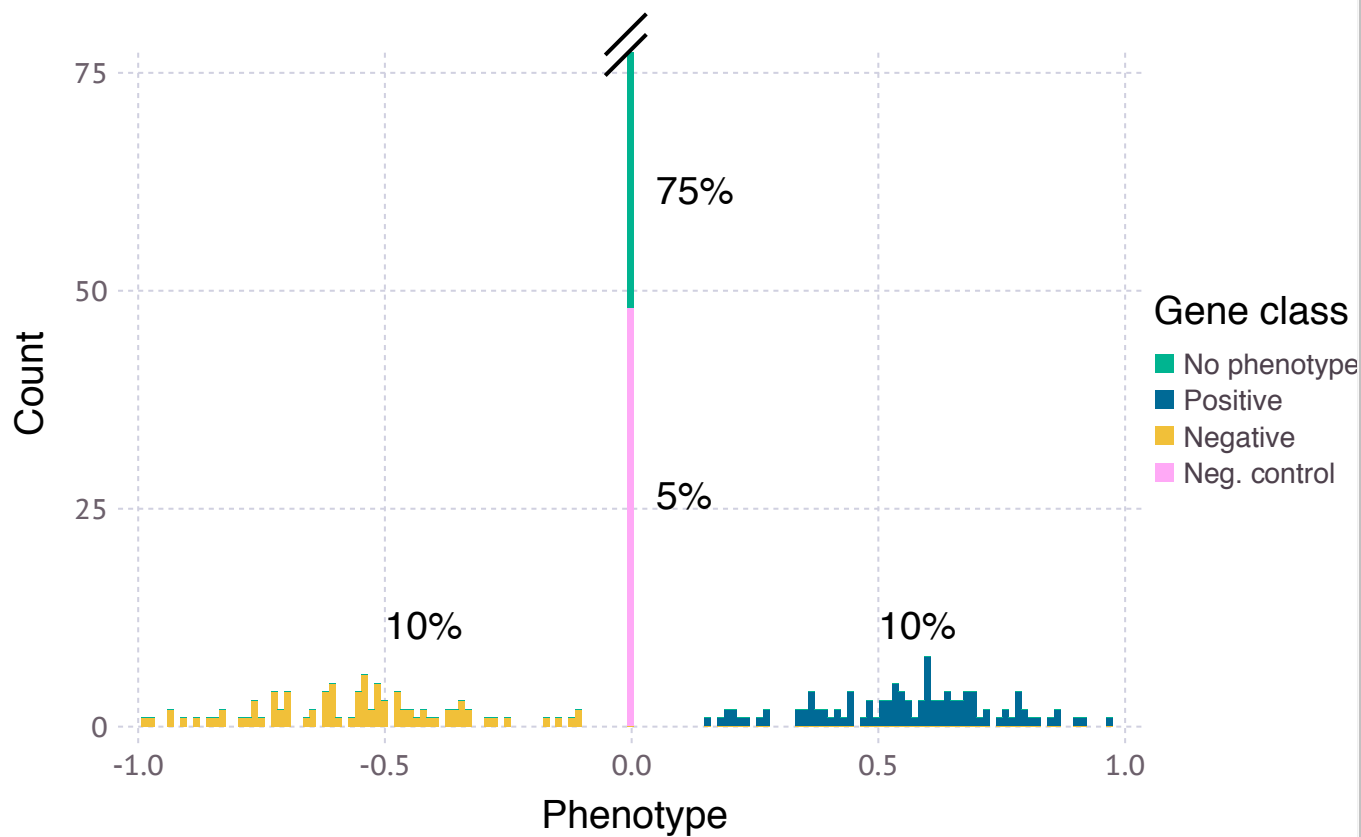
408

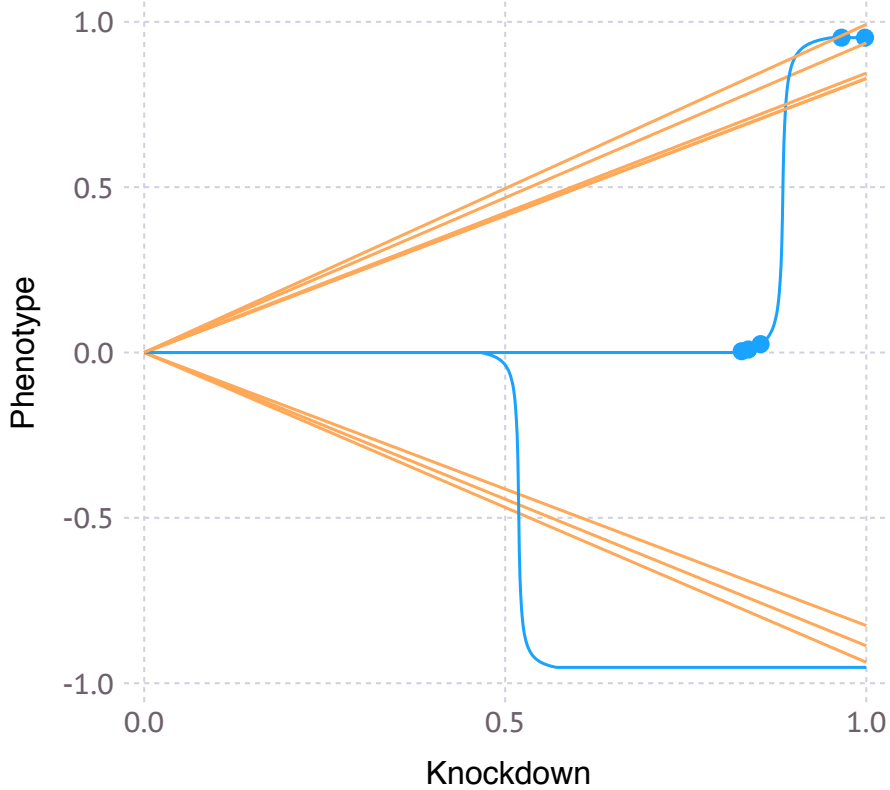
409 **Figure 13.** Comparison of CRISPRn and CRISPRi screen performance for genes with different  
410 knockdown-phenotype relationships. Simulations of FACS-based screens were run for 100x  
411 representation at the transfection, bottleneck and sequencing stages. The simulated genome  
412 contained 75% of genes with a linear knockdown-phenotype relationship and 25% of genes with  
413 a sigmoidal knockdown-phenotype relationship, as defined in the Online Methods. Performance  
414 in hit detection was quantified as AUPRC either for all genes, or only for linear or sigmoidal  
415 genes. Lines and light margins represent means and 99% confidence intervals, respectively, for  
416 100 independent simulation runs.

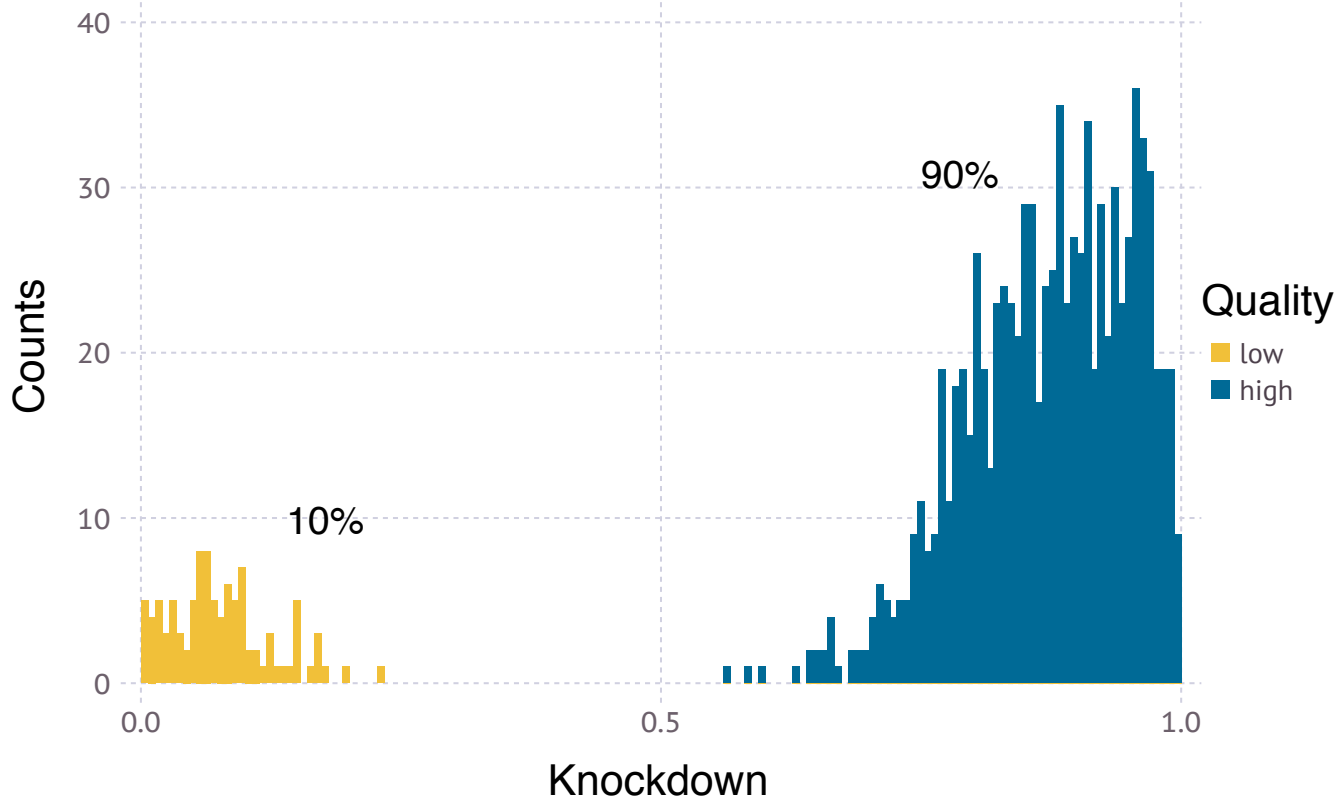
417

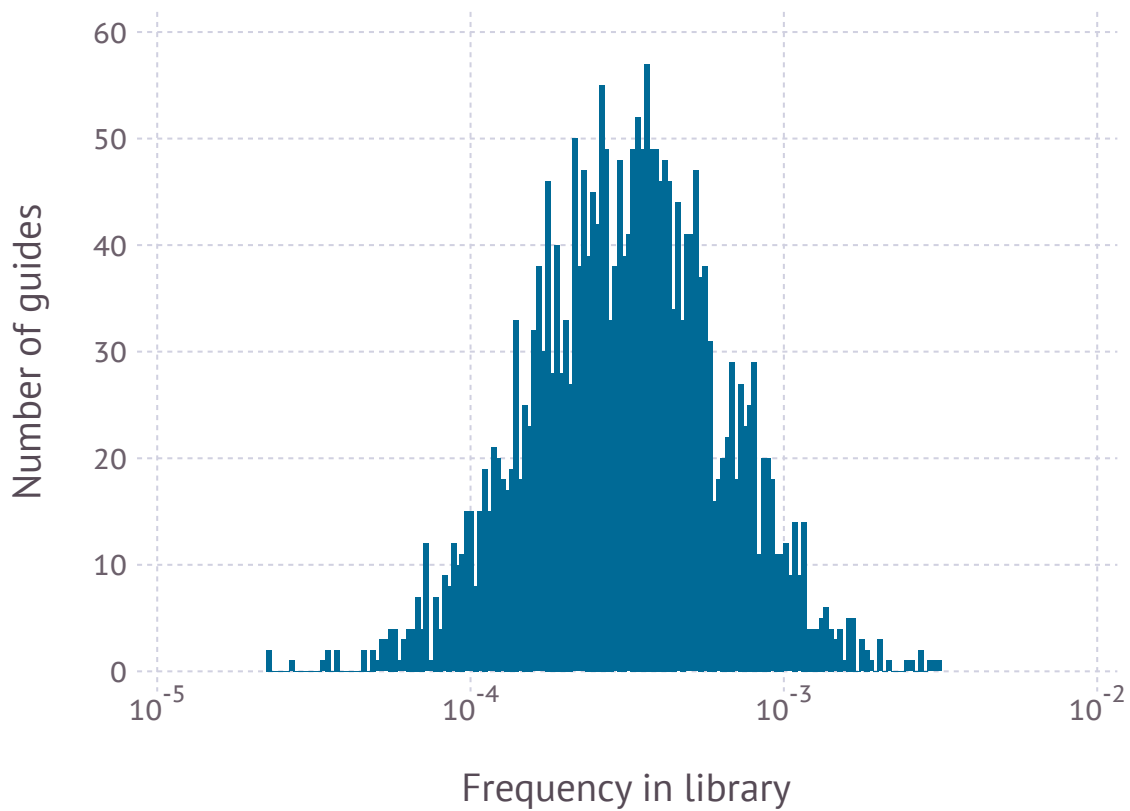
418 **Figure 14.** Experimental data from FACS-based screens resembles simulated data shown in  
419 Figure 6. Grey dots: non-targeting sgRNAs, dots on a red-white-blue color scale: targeting  
420 sgRNAs. Number of deep sequencing reads for each sgRNA in two populations separated based  
421 on a fluorescent reporter signal are shown. **(a)** Screen carried out with high representation at all  
422 stages. **(b)** Screen with low representation at the infection stage. **(c)** Screen with low  
423 representation at the selection stage.





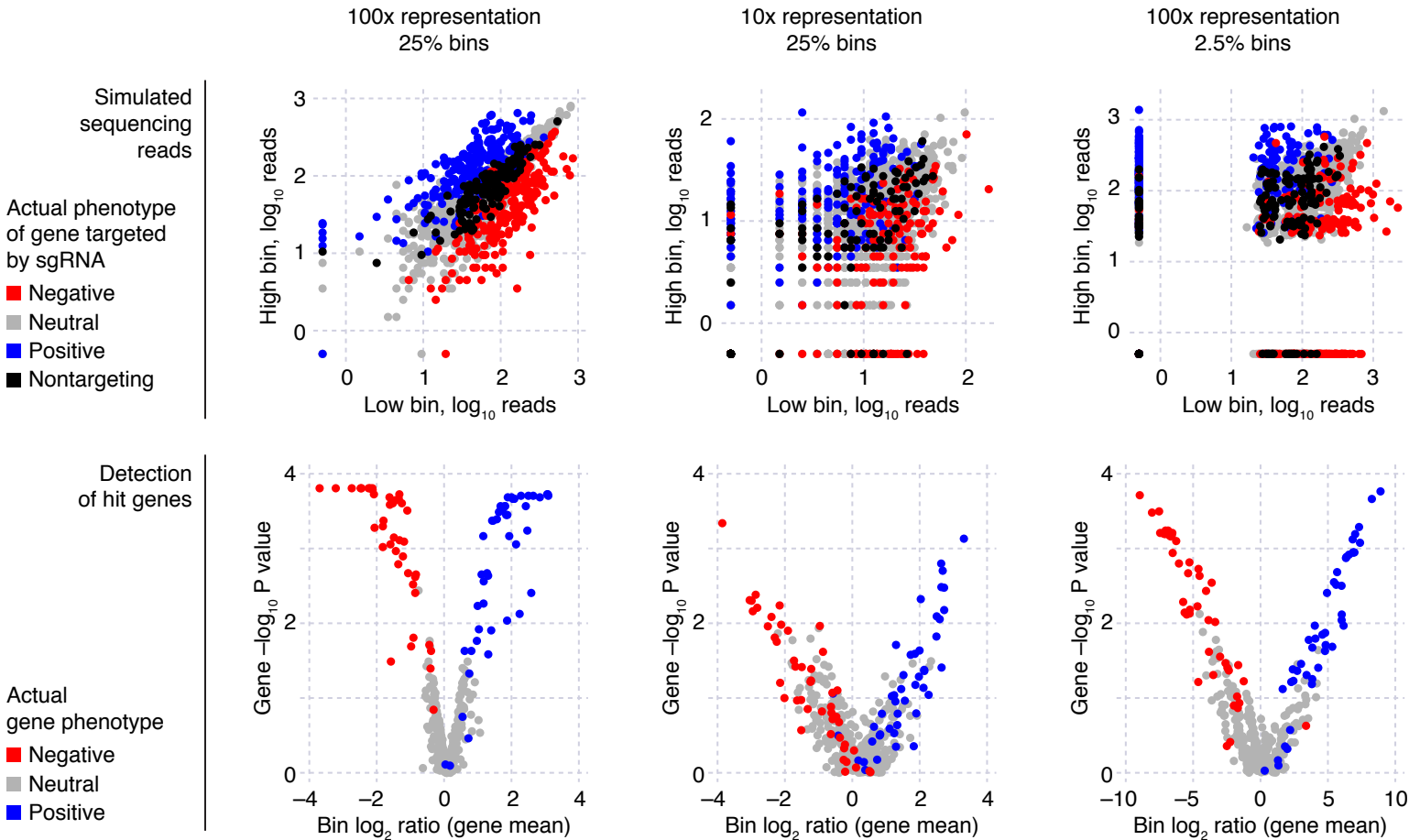




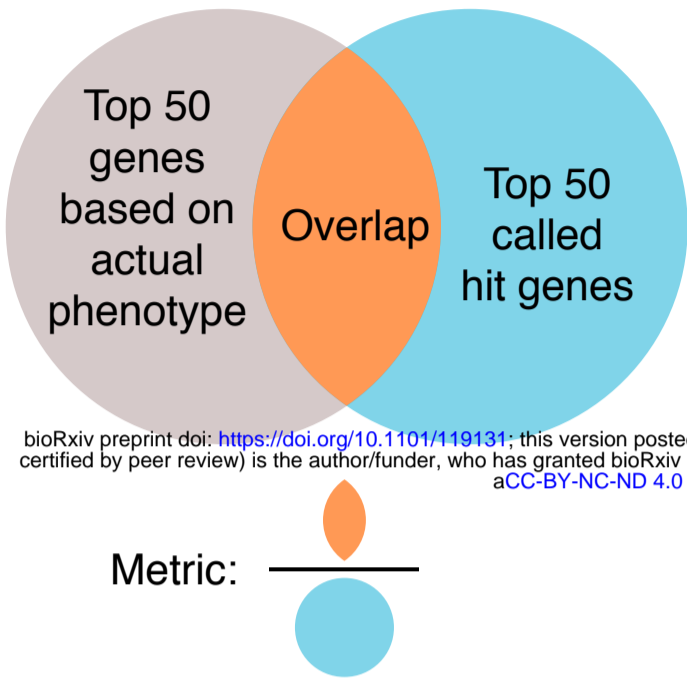




Experimental parameters, FACS-based screen

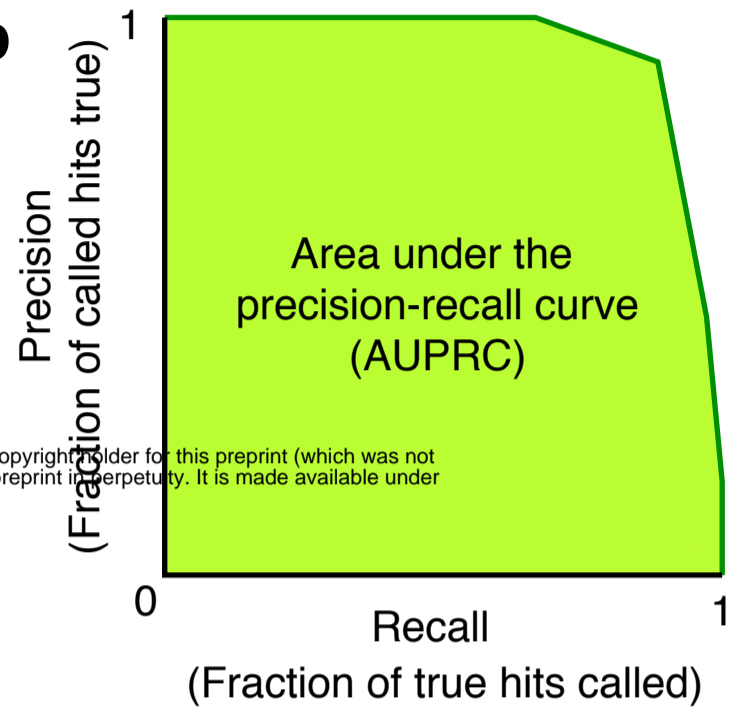


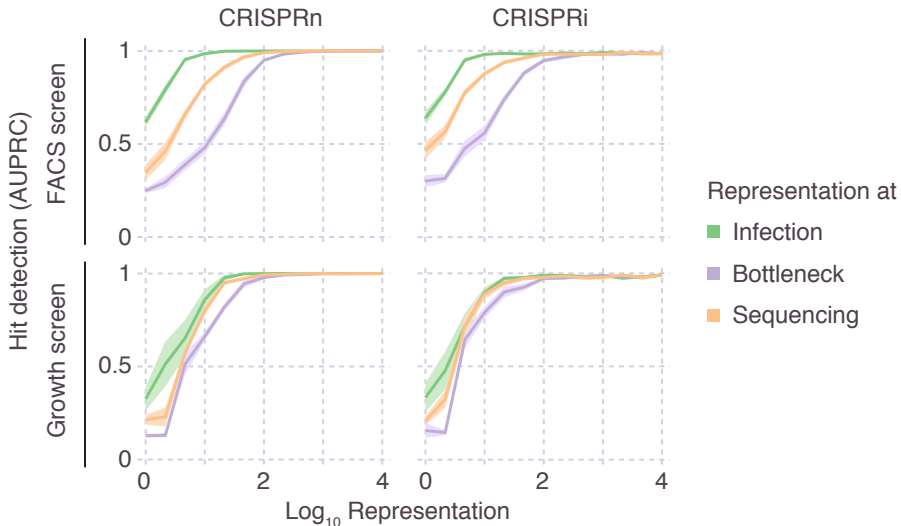
**a**

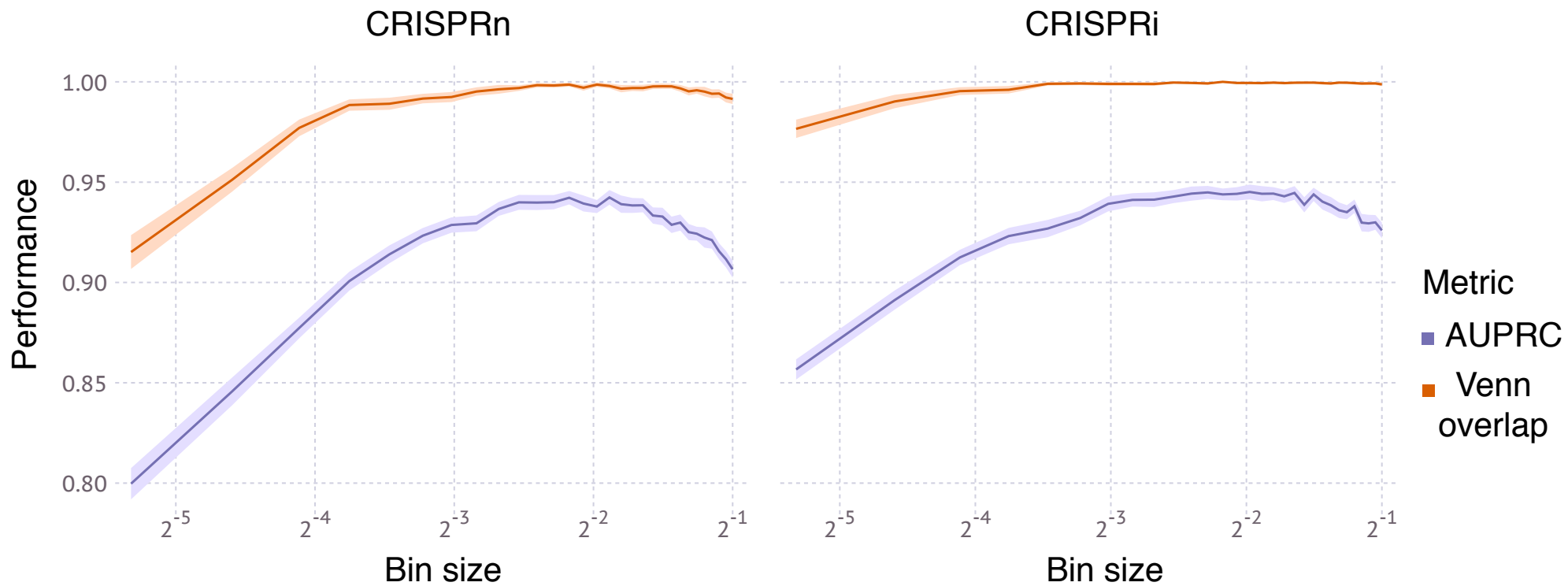


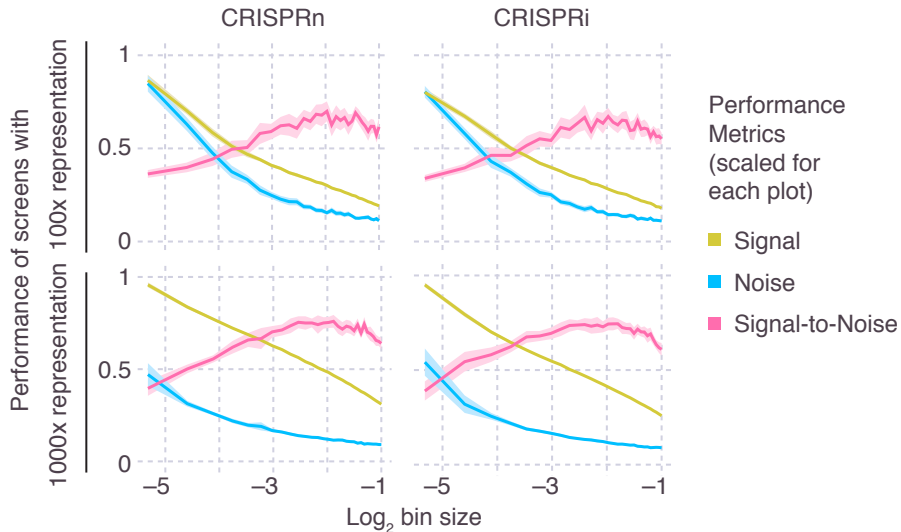
bioRxiv preprint doi: <https://doi.org/10.1101/119131>; this version posted March 22, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

**b**





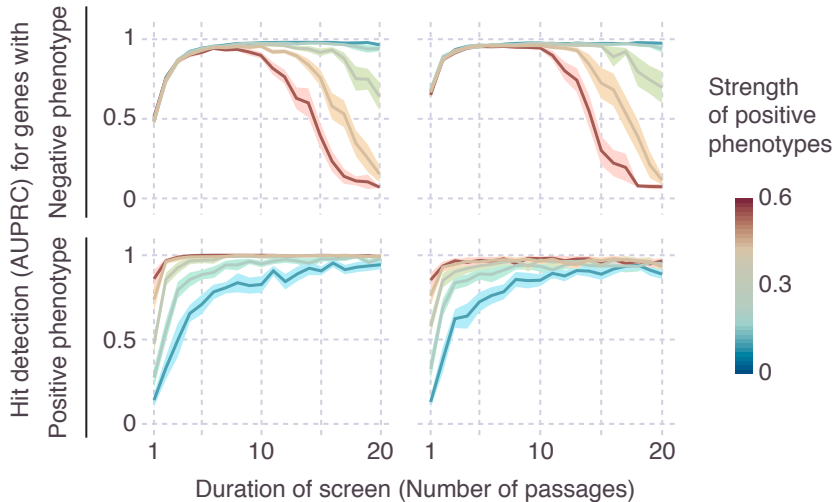




# Growth screen

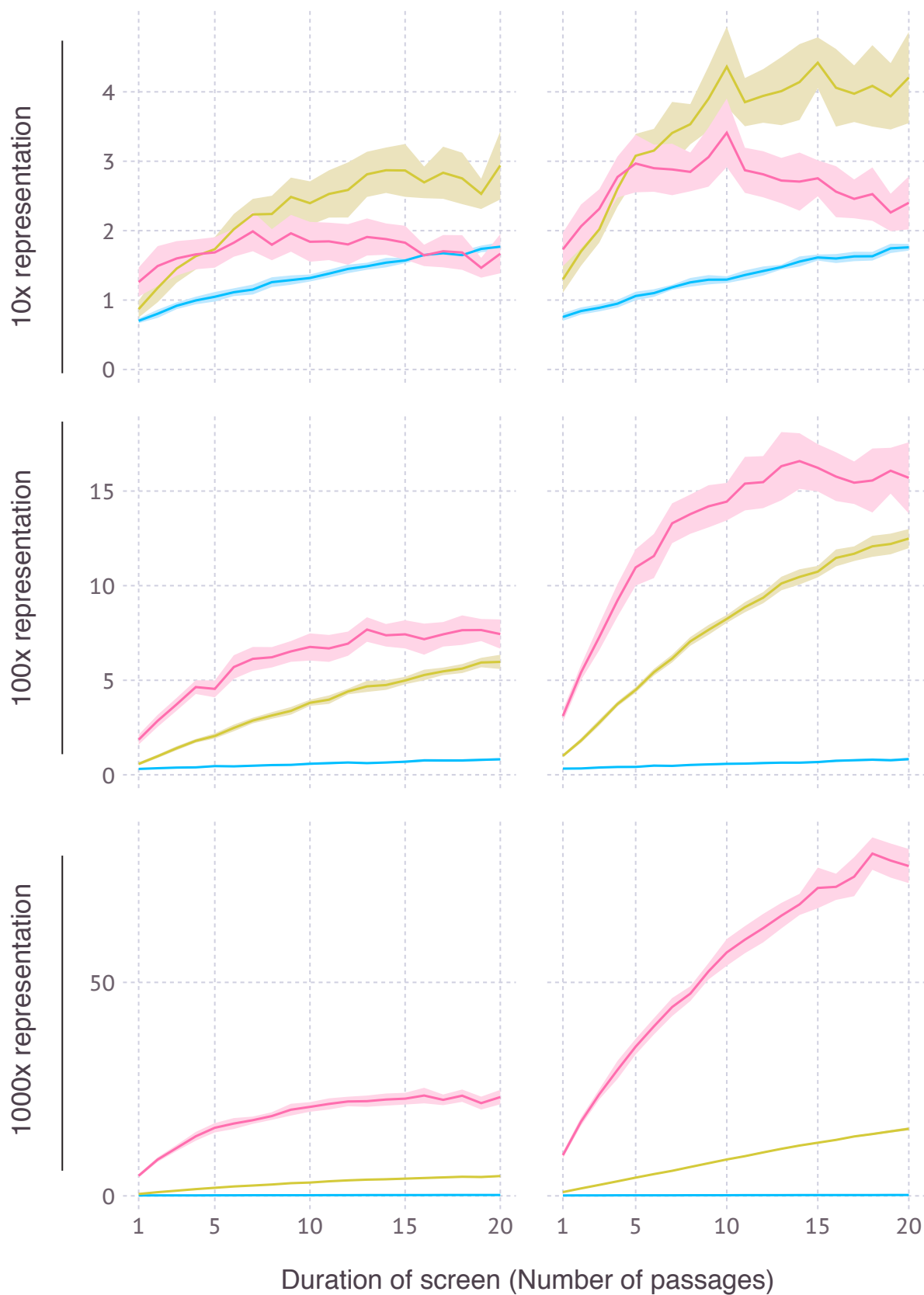
## CRISPRn

## CRISPRi



CRISPRn

CRISPRi



■ Signal

■ Noise

■ Signal-to-Noise

Performance Metrics

