

1

2 **A dose-response model for statistical analysis of chemical genetic interactions in**
3 **CRISPRi screens**

4

5 Sanjeevani Choudhery^{1*}, Michael A. DeJesus², Aarthi Srinivasan¹, Jeremy Rock², Dirk

6 Schnappinger³, Thomas R. Ioerger¹

7

8 ¹Department of Computer Science and Engineering, Texas A&M University, College Station,
9 Texas, United States of America

10

11 ²Laboratory of Host-Pathogen Biology, The Rockefeller University, New York, New York, United
12 States of America

13

14 ³Department of Microbiology and Immunology, Weill Cornell Medical College, New York, New
15 York, United States of America

16

17 * Corresponding author

18 E-mail: schoudhery@tamu.edu (SC)

19

20 Short Title: Statistical Analysis of Chemical Genetic Interactions in CRISPRi data

21

22 Keywords: CRISPRi; chemical genetic interactions; dose-response; drug target identification;

23 statistical analysis

24 **Abstract**

25 An important application of CRISPR interference (CRISPRi) technology is for identifying
26 chemical-genetic interactions (CGIs). Discovery of genes that interact with exposure to
27 antibiotics can yield insights to drug targets and mechanisms of action or resistance. The
28 objective is to identify CRISPRi mutants whose relative abundance is suppressed (or enriched) in
29 the presence of a drug when the target protein is depleted, reflecting synergistic behavior.
30 Different sgRNAs for a given target can induce a wide range of protein depletion and differential
31 effects on growth rate. The effect of sgRNA strength can be partially predicted based on
32 sequence features. However, the actual growth phenotype depends on the sensitivity of cells to
33 depletion of the target protein. For essential genes, sgRNA efficiency can be empirically
34 measured by quantifying effects on growth rate. We observe that the most efficient sgRNAs are
35 not always optimal for detecting synergies with drugs. sgRNA efficiency interacts in a non-linear
36 way with drug sensitivity, producing an effect where the concentration-dependence is
37 maximized for sgRNAs of intermediate strength (and less so for sgRNAs that induce too much or
38 too little target depletion). To capture this interaction, we propose a novel statistical method
39 called CRISPRi-DR (for Dose-Response model) that incorporates both sgRNA efficiencies and
40 drug concentrations in a modified dose-response equation. We use CRISPRi-DR to re-analyze
41 data from a recent CGI experiment in *Mycobacterium tuberculosis* to identify genes that interact
42 with antibiotics. This approach can be generalized to non-CGI datasets, which we show via an
43 CRISPRi dataset for *E. coli* growth on different carbon sources. The performance is competitive
44 with the best of several related analytical methods. However, for noisier datasets, some of these
45 methods generate far more significant interactions, likely including many false positives,

46 whereas CRISPRi-DR maintains higher precision, which we observed in both empirical and
47 simulated data.

48

49 **Author Summary**

50 CRISPRi technology is revolutionizing research in various areas of the life sciences,
51 including microbiology, affording the ability to partially deplete the expression of target proteins
52 in a specific and controlled way. Among the applications of CRISPRi, it can be used to construct
53 large (even genome-wide) libraries of knock-down mutants for profiling antibacterial inhibitors
54 and identifying chemical-genetic interactions (CGIs), which can yield insights on drug targets and
55 mechanisms of action and resistance. The data generated by these experiments (i.e., sgRNA
56 counts from high throughput sequencing) is voluminous and subject to various sources of noise.
57 The goal of statistical analysis of such data is to identify significant CGIs, which are genes whose
58 depletion sensitizes cells to an inhibitor. In this paper, we show how to incorporate both sgRNA
59 efficiency and drug concentration simultaneously in a model (CRISPRi-DR) based on an
60 extension of the classic dose-response (Hill) equation in enzymology. This model has advantages
61 over other analytical methods for CRISPRi, which we show using empirical and simulated data.

62

63

64 Introduction

65 CRISPR technology is becoming an increasingly important tool for genome-wide
66 identification of gene functions in various environmental conditions [1-3]. For example, several
67 different approaches have been devised to exploit CRISPR to induce depletion of target
68 proteins. In the earlier CRISPRko approaches, a nuclease-active form of CAS9 was used to
69 deactivate target genes by cutting the DNA at a target locus and induce DNA repair, which could
70 introduce indels causing frameshifts or inserting novel elements, abrogating their function
71 completely [1-3]. Another approach, CRISPRa, utilizes dCAS9 fusions with effectors that actively
72 enhance or suppress transcription through direct interaction with the RNA polymerase (such as
73 transcription factors that can activate transcription) [4].

74 In CRISPR interference (CRISPRi), a catalytically-dead CAS9 protein (dCAS9) is recruited to
75 a chromosomal locus by a single guide RNA (sgRNA) with a short (~20 bp) complimentary
76 sequence and physically blocks transcription [5]. dCAS9 nucleases from several different
77 organisms are available for CRISPRi (e.g. *S. pyogenes*, *S. thermophilus*, [6]) and different
78 promoters and chemicals have been used for dCAS9 induction. The degree of CRISPR
79 interference can be tuned by modulating the level of dCAS9 expression [7], varying the sgRNA
80 sequence with respect to its length, GC-content, targeting sequence complementarity, position
81 in the gene, or similarity of targeted PAM (protospacer adjacent motif) sequence, to consensus
82 for optimal dCAS9 recognition, [5, 6, 8-11]. While in mammalian systems, efficiency of sgRNAs
83 can vary among multiple cell types, [9], for simplicity, our focus is on studying single defined
84 lineages, as in bacterial strains. Tuning CRISPRi allows to deplete the targeted gene product to
85 intermediate levels [5], which allowed the introduction of the concept of gene ‘vulnerability’ as

86 describing the sensitivity of cells to partial depletion of individual proteins [12]. By this
87 definition, highly vulnerable genes are genes for which even small depletion of the encoded
88 protein causes growth impairment, which can be quantified efficiently on a genome-wide scale
89 using high-throughput sequencing [12]. The vulnerability of a gene can be both condition
90 dependent and strain or cell type dependent [12].

91 One interesting application of CRISPRi is to reveal targets of antibiotics or mechanisms of
92 resistance through chemical-genetic interactions (CGI) [7, 13]. CRISPRi libraries can be designed
93 to contain multiple sgRNAs targeting each gene, resulting in a set of thousands of individual
94 depletion mutants [12]. In this context, ‘mutant’ refers to a cell line transformed with a
95 integrative plasmid capable of expressing the dCAS9 protein and the unique targeting sgRNA,
96 even though it contains the wild-type gene sequence. The abundance of each mutant can be
97 quantified by amplifying the sgRNA targeting sequence which functions as a molecular barcode,
98 and then performing deep sequencing to count the number of barcodes for each sgRNA in a
99 treatment [6]. The analysis of such datasets is challenging, due to various sources of noise
100 which introduce variability in the counts.

101 There are several previously published methods for statistical analysis of CRISPR
102 datasets. One, called MAGeCK [14] (originally intended for CRISPRko screens), calculates a log-
103 fold-change (of mean counts) for each sgRNA between a treatment condition and a reference
104 condition (control), and uses a Gaussian distribution to estimate the significance of differences
105 in mean sgRNA abundance between treatments and controls (based on the implementation in
106 the source code, which differs from the description in the publication). To evaluate effects at
107 the gene level, individual sgRNAs are combined in MAGeCK using Robust Rank Aggregation

108 (RRA) to prioritize genes whose sgRNAs show greater enrichment or depletion on average than
109 other genes in the genome. MAGeCK has been used for evaluating chemical-genetic interactions
110 (CGI) with antibiotics [14]. A variant called MAGeCK-MLE [15] fits a Bayesian model by
111 Maximum Likelihood that captures changes in mean counts with increasing time or
112 concentration, along with effectiveness of each sgRNA through posterior probabilities of a
113 binary variable, to determine the overall probability that a gene interacts. Other approaches
114 such as CRISPhieRmix [16] use mixture models to separate effective from ineffective sgRNAs,
115 and thereby identify interacting genes as those containing a significant subset of effective
116 sgRNAs. DrugZ [17] identifies significant interactions by averaging together Z-scores (assuming a
117 Normal distribution) of log-fold-changes of sgRNAs at the gene level. DEBRA [18] utilizes
118 DeSeq, a method for transcriptomic analysis, which employs the Negative Binomial distribution
119 for counts and a more sophisticated method for modeling variance and using it to discriminate
120 genes displaying significant changes in mean counts.

121 However, most of these methods have one of two limitations when applied to identify
122 genes affecting drug potency. First, CGI experiments are ideally carried out with multiple drug
123 concentrations around the MIC (minimum-inhibitory concentration), since it is often difficult to
124 anticipate what concentration will stimulate the right amount of growth inhibition in
125 combination with CRISPRi-induced depletion of target proteins. However, many of the existing
126 methods analyze the data for each drug concentration independently (i.e. comparing each
127 concentration to a no-drug control). Since knock-down mutants might exhibit depletion at one
128 concentration but not others, results from multiple concentrations must be combined post-hoc.
129 As an example, the authors in [13] chose to combine results from analyzing different

130 concentrations of a given drug using MAGeCK-RRA by taking the union of significant interacting
131 genes at each individual concentration. Due to the noise in these CRISPRi experiments,
132 analyzing concentrations independently increases the risk of detecting false positives (in the
133 sense that non-interacting genes might be spuriously called as hits at different concentrations).

134 Second, many of the analytical methods do not explicitly take into account differences in
135 sgRNA efficiency (i.e. take sgRNA efficiencies as an input in the model). Different sgRNAs can
136 induce different degrees of depletion of their target genes, and this in turn causes different
137 effects on growth rate, depending on sensitivity of the cells to protein depletion [10]. This can
138 be quantified beforehand by evaluating the growth rate of individual CRISPRi mutants (with
139 unique sgRNAs) in a growth experiment and determining the actual fitness defect caused by
140 target knockdown [11, 12]. In highly vulnerable genes, the effect of protein depletion by sgRNAs
141 on cell growth rate (efficiency) can span a range from no effect to severe growth defect. Early
142 applications of CRISPR were primarily being used to fully inactivate genes (e.g. CRISPRko), rather
143 than to produce graded depletion effects. Therefore, at the time some of these methods were
144 developed, this information was often not used, as methods to quantify sgRNA efficiencies were
145 not well developed. Even in MAGeCK, the Robust Rank Aggregation method treats all sgRNAs
146 in a gene as “equal” a priori, without differentiating them based on the expected effects due to
147 sgRNA efficiency. (Efficiency is not an input.) In contrast, it has been recognized that different
148 sgRNAs can have different efficiency, and several papers have investigated the factors that are
149 associated with stronger sgRNAs [19], especially sequence-based attributes such as similarity to
150 optimal PAM sequence, length and GC content of targeting sequence, mismatches, etc. [5, 8,
151 10]. Mathis, Otto and Reynolds (11) exploit this to synthetically create a diverse set of sgRNAs

152 with a range of efficiencies by mutating the guide RNA sequences, which they quantify by
153 empirically fitting growth curves for each modified sgRNA with a logistic equation. Interacting
154 genes are then found using differences in the fitted parameters that includes the quantified
155 growth rates and the Hill coefficient. Among all the existing CRISPR analytical methods,
156 MAGeCK-MLE [15] is the only other method that explicitly includes sgRNA efficiencies as an
157 input, which are used to set the prior probabilities that each sgRNA is effective or not (because
158 of their focus on CRISPRko) in the joint probability formula, to initialize for the Expectation
159 Maximization iterations.

160 In the application to CGI data, a regression model can be used to integrate data over
161 multiple drug concentrations [20]. The degree of a gene-drug interaction is reflected by the
162 coefficient (or slope) for the dependence of CRISPRi mutant abundance on drug concentration.
163 This regression approach was previously introduced in CGA-LMM for analysis of hypomorph
164 libraries (where there is typically just one mutant representing each gene) [20]. It was based on
165 the theory that depletion of the target of a drug should ideally synergize with increasing
166 concentrations of the drug. While exposure to an inhibitory compound will challenge the
167 growth of all the mutants in a hypomorph library, mutants with depletion of a gene that
168 interacts with a drug (e.g. prototypically, an essential gene that is the drug target) will exhibit
169 excess depletion relative to others in the library due to the combined effect of both the growth-
170 inhibition due to the drug treatment in conjunction with the growth-impairment due to knock-
171 down of an vulnerable gene, making these hypomorphic mutants even more sensitive to the
172 drug. For genes that genuinely interact with a given drug, this depletion effect should be
173 exacerbated at higher drug concentrations (i.e. be dose-dependent); thus, genes of greatest

174 relevance would be those that exhibit concentration-dependent effects. While the (log of)
175 abundance of a depletion mutant does not have to decrease perfectly linearly with the (log of)
176 drug concentration to obtain a significant negative coefficient (slope) in the regression, there
177 should be a general trend supporting that relative abundance decreases as concentration
178 increases.

179 One of the challenges in extending this prior regression approach (CGA-LMM) to CRISPRi
180 screens was incorporating information on sgRNA efficiency. Even in essential genes, some
181 sgRNAs may produce strong depletion of the target, while others might be almost completely
182 ineffective. While sgRNA strength can be partially predicted (with intermediate accuracy) from
183 sequence alone [9, 12], the actual growth phenotype depends on vulnerability of the target
184 gene (sensitivity of cells to depletion of the protein product), which is what is meant by sgRNA
185 efficiency. Even sgRNAs that are predicted to be strong might not cause a growth defect if they
186 are in a non-essential gene. sgRNA efficiency must be empirically quantified by measuring
187 growth rates in standard growth media (e.g. by fitting exponential growth curves based on
188 optical density, or using a reporter gene) with versus without induction of dCAS9, and then
189 calculating relative fitness defects [11]. An alternative approach is to fit the abundance of
190 depletion mutants to a piecewise linear model that allows for a preliminary lag phase, and then
191 extrapolating the model to predicted log-fold-change (LFC) at a fixed number of generations
192 [12]. Any such measure of sgRNA efficiency can be incorporated as a term in the CRISPRi-DR
193 model we present below. Although one could contemplate adding the efficiency of each sgRNA
194 into a simple regression model to predict abundances for each gene, a significant problem

195 (expanded upon below) is that sgRNAs of different efficiency can show different concentration
196 dependence, resulting in non-linear interactions among variables.

197 In this paper, we propose a modified regression approach for CRISRPi data (called
198 CRISPRi-DR) that incorporates both drug concentration and sgRNA efficiency. The approach is
199 based on the classic dose-response (DR) model for inhibition activity of drugs; the activity of a
200 target protein typically transitions from high to low in the shape of an S-curve as concentration
201 increases (on a log scale), which can be modeled with a Hill equation. The parameters of the Hill
202 equation for a given drug can be fit by performing a log-sigmoid transformation of the mutant
203 abundance data and then using ordinary least-squares regression. We show how sgRNA
204 efficiency can be incorporated into this model as a multiplicative term in the Hill equation,
205 which becomes an additive effect in the log-sigmoid transformed data. The benefit of this
206 model is that it decouples the concentration-dependence from the sgRNA efficiency, so they
207 can be fit as independent (non-interacting) terms in the regression, which ultimately amplifies
208 effects that may be apparent only for a subset of sgRNAs in an optimal efficiency range.

209 CRISPRi-DR is applicable to libraries where there are multiple sgRNAs representing each
210 gene with a range of efficiencies, which can be quantified empirically as an effect on growth rate
211 (fitness defect). The diversity of efficiencies is useful for identifying synergistic effects with
212 treatments/conditions. Thus, the main requirements for CRISPRi-DR are that: a) there are
213 multiple sgRNAs for each target in the library, b) the sgRNAs vary in predicted strength, and c)
214 the actual efficiencies of the sgRNAs (i.e. growth defects due to target depletion) have been
215 experimentally quantified in control conditions, as an input to the analysis method. The
216 primary use case we focus on is identification of chemical-genetic interactions, with drug

217 concentration as a covariate. We demonstrate the value of the CRISPRi-DR analysis method by
218 re-analyzing the data from a recent paper using CRISPRi for chemical-genetic interactions to
219 identify targets of antibiotics in *M. tuberculosis*. However, the approach can be generalized to
220 analyze experiments with other covariates, such as time-points of a treatment, where there is a
221 sigmoidal response in growth. We illustrate this by using CRISPRi-DR to analyze an *E. coli*
222 CRISPRi dataset from an experiment to determine genes differentially required for growth on
223 different carbon sources [11].

224

225

226 **Methods**

227 The CRISPRi-DR method applies to CRISPRi experiments that involve using high-
228 throughput sequencing to tabulate sgRNA counts representing abundance of individual CRISPRi
229 mutants in a population (pooled culture). Each mutant has an sgRNA (on a plasmid) mapping to
230 a target gene that can reduce its expression (e.g. with dCAS9 induction). In CGI applications, the
231 culture is treated with antibiotics or inhibitors at various concentrations, along with a no-drug
232 control, and DNA is extracted, PCR-amplified, and sequenced, producing counts representing
233 each sgRNA. If Y_{ijk} is the abundance (i.e. count) for an sgRNA i in a condition j for replicate k ,
234 normalized abundance can be calculated by $Y'_{ijk} = \frac{Y_{ijk}}{\sum_{x=1}^n Y_{xjk}}$, where each count is divided by the
235 sum of counts of all the sgRNAs observed in a given condition and replicate. Let U'_i be the
236 normalized abundance of sgRNA i in the uninduced condition, then the normalized relative

237 abundances of an sgRNA i in all induced samples can be calculated as: $A_{ijk} = \frac{Y'_{ijk}}{U'_{ri}}$, assuming
238 that the counts in the uninduced condition represents full abundance of each clone (normal
239 growth without target depletion).

240

241 CRISPRi Dose-Response Model

242 The CRISPRi-DR model for analyzing CRISPRi data from CGI experiments is an extension
243 of the basic dose-response model, extended to incorporate sgRNA efficiencies. The dose-
244 response effect of an inhibitor on the activity of an enzyme is traditionally modeled with the
245 Hill-Langmuir equation.

$$246 \quad \theta = \frac{1}{1 + \left(\frac{K_A}{[L]}\right)^n} \quad (1)$$

247 where θ is the fraction of abundance (relative to no drug), $[L]$ is the ligand concentration, K_A is
248 the concentration at which there is 50% activity and n is the Hill coefficient.

249 Applying Eq (1) to the CGI data, the relative abundance of sgRNAs A_{ijk} is used as the
250 predictor variable and $[D_j]$ is the concentration of drug j that the k th replicate count of sgRNA i
251 was extracted from,

$$252 \quad A_{ijk} = \frac{1}{1 + \left(\frac{IC_{50}(D_j)}{[D_j]}\right)^{H_d}} \quad (2)$$

253 The unknown parameters are the IC_{50} value (inhibitory concentration that causes 50% growth
254 inhibition) and the Hill coefficient H_d . The plot of the concentration versus relative abundance

255 of an sgRNA (A_{ijk}) produces a sigmoidal curve, demonstrating how activity decreases as
256 concentration increases, with the IC_{50} , representing the mid-point of the transition.

257 The dose-response model seen in Eq 2 can be extended to account for sgRNA efficiency
258 by incorporating a multiplicative factor in the denominator:

$$259 \quad A_{ijk} = \frac{1}{1 + \left(\frac{IC_{50}(D_j)}{[D_j]} \right)^{H_d} \left(\frac{K_s}{S_i} \right)^{H_s}} \quad (3)$$

260 sgRNA efficiency, S_i , is an empirical measure of the degree of growth impairment resulting from
261 target depletion. This can be assessed in several ways, such as estimating change in exponential
262 growth rate in a reference condition in a growth experiment [21]. Alternatively, Bosch et al [12]
263 use estimated log-fold change of abundance (induced vs uninduced) at a fixed number of
264 generations of growth in-vitro in the absence of drug, extrapolated from a model fit to empirical
265 data (passaging experiment) that allows for a lag phase. K_s represents the unknown
266 intermediate sgRNA efficiency that causes 50% depletion of mutant abundance (half-way
267 between no depletion and full depletion), and the H_s is the unknown Hill coefficient that
268 represents how sensitive mutant abundance is to depletion of the target protein.

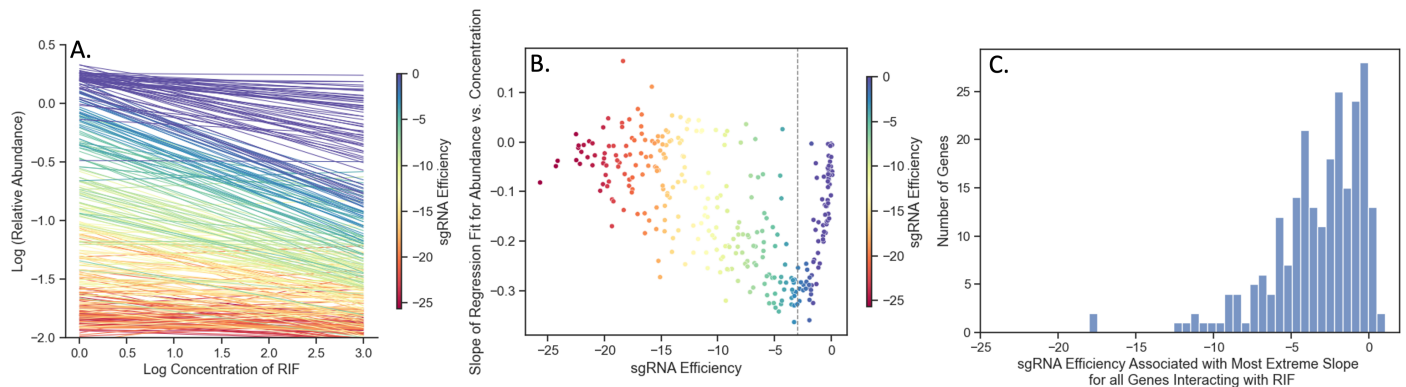
269

270 **Relationship between drug concentration and gene depletion within** 271 **the CRISPRi-DR model**

272 Abundance of mutants in a CRISPRi CGI experiment can be affected simultaneously by
273 both presence of an inhibitor and depletion of an interacting gene. However, the concentration-
274 dependent effect of a drug on mutant abundance can be different for sgRNAs of different

275 efficiency. Fig 1 illustrates the interaction between these two effects for *rpoB* (RNA polymerase
276 beta chain) in an *Mtb* CRISPRi library treated with rifampicin with 5 days pre-depletion. The
277 lines in Fig 1A are regression fits obtained for each sgRNA in *rpoB* using regression of log
278 abundances against log concentration of rifampicin, $\log(A_{ijk}) = C + B \cdot \log([D_j])$, where C is
279 in the intercept and B is the slope of the regression, representing concentration dependence,
280 and $\log(A_{ijk})$ are log relative abundances obtained as described above. The left-most side of
281 Fig 1A shows the range of abundances in the no-drug control (induced library in media without
282 rifampicin). These differences in abundances (dispersion along Y-axis) are due solely to the
283 growth impairment caused by depleting RpoB. As concentration of RIF increases, some of the
284 sgRNAs show very negative slopes, while other sgRNAs show slopes closer to 0. A parabolic-
285 type curve emerges in Fig 1B when the slopes B from the regressions are plotted against the
286 sgRNA efficiencies. Both the most efficient sgRNAs (colored red) and the least efficient sgRNAs
287 (purple) have slopes around 0 (no concentration dependence). Highly efficient sgRNAs (red) can
288 cause excessive depletion (even without drug), making it difficult to detect additional decreases
289 due to increasing drug concentration. Comparatively, sgRNAs with very low efficiency (purple)
290 might not induce enough depletion to synergize with the drug. The sgRNAs surrounding the
291 minimum point of the parabolic curve (dashed line) in Fig 1B reflect those of intermediate
292 efficiency where the ability to detect synergy with the drug is maximized. These are the sgRNAs
293 in Fig 1B that show the most negative slope with increasing concentration (dark green-indigo).
294 As Fig 1C shows, the efficiency where the slopes reach their extremes (most negative; or most
295 positive for those showing enrichment) can be different for each gene but tend to fall in an
296 intermediate region of sgRNA efficiency (0 to -5). The histogram shows that sgRNA efficiency at

297 which the most extreme (largest or smallest) concentration-dependent slope is achieved over
298 all interacting genes (236 for RIF D5). Hence, the sgRNAs that are optimal for detecting CGIs are
299 not necessarily the strongest (most efficient). The variability of concentration-dependence
300 (slope) with sgRNA efficiency suggests a possible non-linear interaction between the variables.
301 This nonlinearity is captured in the multiplicative terms of the dose-response model (Eq (3)).



302

303 **Fig 1. Effect of sgRNA efficiency on concentration dependence for sgRNAs in *rpoB* in a**
304 **CRISPRi library treated with RIF (D5).**

305 (A) Regression lines for log(relative abundance) against log(concentration) for all sgRNAs
306 in *rpoB* in a library treated with RIF for 5 days pre-depletion. The lines that reflect the
307 extremes of the sgRNA efficiency (red or purple), are flat and do not show much change
308 in abundance. Comparatively, intermediate sgRNA efficiency (dark green to indigo)
309 shows the most negative slopes, reflecting maximum synergy with drug. (B) Comparison
310 of sgRNA efficiency and slopes of the regressions seen in Panel A for each sgRNA. Each
311 point is an sgRNA colored by its efficiency. The most efficient sgRNAs (purple) and the
312 least efficient sgRNAs (red) show concentration slopes around 0. The dotted line reflects
313 the minimum of the parabolic curve. (C) Histogram of sgRNA efficiencies where the
314 slopes reach their most extreme (positive or negative) for 236 interacting genes in RIF

315 D5. The distribution shows that most of the extrema sgRNAs for interacting genes fall in
316 the range of -5 to 0 (note: not the strongest sgRNAs, which would have efficiencies
317 around -25).

318

319 **Linearization and parameter estimation**

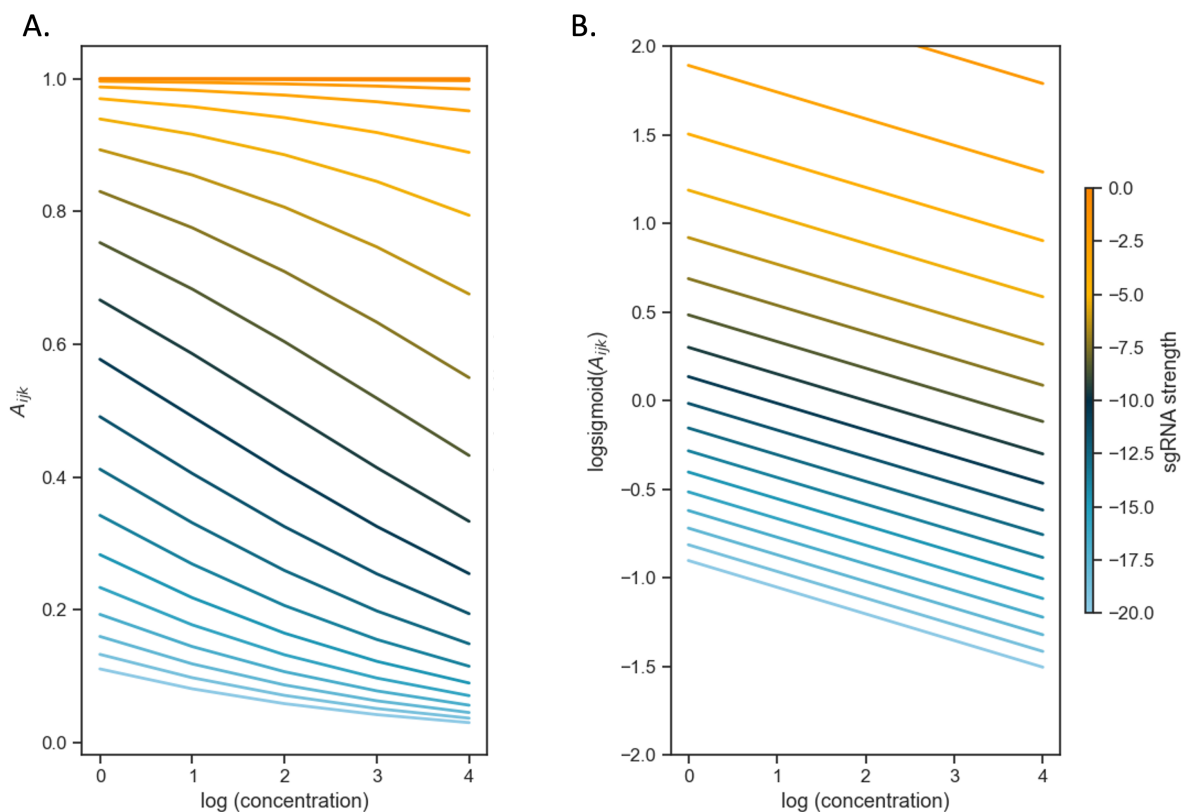
320 The dose-response model Eq (3) can be linearized through a log-sigmoid transformation.

$$321 \quad \text{Log} \left(\frac{A_{ijk}}{1 - A_{ijk}} \right) = H_d \cdot \log([D_j]) + H_s \cdot S_i + C$$

$$322 \quad C = H_s \cdot \log(K_s) - H_d \cdot \log(IC_{50}(D_j)) \quad (4)$$

323 In this log-sigmoid transformed space, the concentration-dependence and effect of sgRNA
324 efficiency have been decoupled, appearing as independent linear terms with the Hill coefficients
325 (H_s and H_d) as the variables to solve for by a standard regression. The inflection parameters of
326 the sigmoid curve (K_s and IC_{50}) are combined in the intercept C in the model. Importantly, this
327 model implies that the effects of growth impairment due to the depletion of a vulnerable gene
328 and growth inhibition due to the drug on the overall (relative) abundance of a given mutant
329 become additive in this log-sigmoid-transformed space. To illustrate this, the CRISPRi-DR
330 equation is simulated by plotting idealized relative abundances (in Fig 2) using parameters
331 chosen to emulate what is seen in Fig 1A, the plot of slopes over a systematic range of sgRNA
332 efficiencies and drug concentrations for *rpoB*. In Fig 2A, the slopes of the concentrations are
333 plotted against abundances calculated using the dose-response model. The slopes vary as a
334 function of the starting depletion (left-hand side), which is due to sgRNA efficiency alone
335 (colored gradient based on sgRNA efficiency value). The slopes are most negative for

336 intermediate sgRNA efficiency, colored with a dark blue-green hue representing sgRNA
337 efficiency around -10. Fig 2B illustrates the result of the linearization (log-sigmoid
338 transformation) of the Hill equation. All the individual sgRNA regression lines over concentration
339 become parallel, eliminating the dependence on sgRNA efficiency, and allowing them to be fit
340 by a single common slope representing the concentration-dependence averaged over all the
341 sgRNAs.



342
343 **Fig 2. The log-sigmoid transformation of abundances allows the CRISPRi-DR model to**
344 **factor in the non-linear effect of sgRNA strength on concentration dependence. (A)**
345 Simulation of sgRNAs abundances for an ideal essential gene. Parameters used in
346 simulation: $H_s = -4$, $IC_{50} = 0.25$, $K_s = -10$ and $H_d = -0.5$ over a range of sgRNA efficiencies
347 and drug concentrations. (B) When the log-sigmoid transformation of the abundances is

348 applied, we see all the regression fits are parallel to one another, allowing to be fit by a
349 single common slope, representing the concentration dependence over all sgRNAs,
350 regardless of sgRNA efficiency.

351

352 Experimental data (i.e. counts from sequencing, converted to relative abundances for
353 mutants with each sgRNA) are fit on a gene-by-gene basis using ordinary least-square (OLS)
354 regression by the following formula:

$$355 \quad \log\left(\frac{A_{ijk}}{1-A_{ijk}}\right) = \beta_0 + \beta_c \cdot \log([D_j]) + \beta_s \cdot S_i \quad (5)$$

356 where A (relative abundance for each CRISPRi mutant at given drug concentration), S_i (sgRNA
357 efficiency) and $[D_j]$ (concentration of drugs) are columns of a melted matrix. To include the
358 control samples (no-drug, dCAS9-induced controls) in the regression, they are treated as one
359 two-fold dilution lower than the lowest available concentration tested for the drug (to avoid
360 taking the log of 0). Since the log-sigmoid transform of the relative abundances is taken, they
361 must be within the range of (0,1). Although relative abundances greater than 1.0 are possible in
362 treated conditions (relative to uninduced, no-drug controls), especially in cases where target
363 depletion confers a growth advantage and consequent enrichment, we use a squashing function
364 to ensure the relative abundances range between 0 and 1, which is required to take the log-
365 sigmoid transform.

$$366 \quad A_{ijk} = \tau + \frac{(1 - \tau)(1 - e^{-2A_{ijk}})}{(1 + e^{-2A_{ijk}})} \quad (6)$$

367 where $\tau=0.01$ is a pseudo count needed to make abundances non-zero for taking logarithms.

368 Relative abundances that are greater than 1.0 are mapped to just below 1.0, though the

369 mapping is monotonic, so the order among sgRNAs is still preserved (higher abundances
370 become exponentially closer to 1.0).

371

372 **Significance Testing**

373 The statistic that indicates the degree of interaction of each gene with a given drug is the
374 coefficient for the $\log([D])$ term (i.e. slope) in the model. To determine whether the interaction
375 is statistically significant, a Wald test [22] is applied to calculate a P-value reflecting whether the
376 coefficient is significantly different than 0, adjusting for a target FDR (false discovery rate) of 5%
377 over the whole genome using the Benjamini-Hochberg procedure [23]. However, the Wald test
378 by itself yields many genes predicted to interact with the drug (often thousands) with adjusted
379 P-value < 0.05. The test selects genes with slopes that are technically different than 0, but not
380 necessarily large enough to be relevant to the drug mechanism. Our assumption is that most of
381 genes in the genome do not interact with a given drug (at least not directly involved in the
382 mechanism of action or resistance). Many genes have small positive and negative slopes,
383 possibly due to some source of noise in the experiment or generalized phenotypic interactions,
384 which should be filtered out. Therefore, genes are filtered based on the magnitude of the slopes
385 (analogous to the requirement of $|LFC| > 1$ used by Li, Poulton (13) to filter significant genes by
386 MAGeCK). The distribution of slopes over all genes is assumed to be a Normal distribution, and
387 Z-scores are computed for every gene g : $Z_g = \frac{\beta_{c,g} - \mu(\beta_c)}{\sigma(\beta_c)}$, where $\sigma(\beta_c)$ is the standard
388 deviation of the slopes of log concentration dependence and $\mu(\beta_c)$ is the mean of the slopes.
389 Genes with $|Z_g| < 2.0$ are filtered out. This produces hits whose slopes are significant outliers

390 ($>2\sigma$) from the rest of the population (i.e. genes in the genome). There are two groups of hits,
391 corresponding to the two tails of the distribution: enriched hits where $Z_g > 2.0$, and depleted
392 hits, $Z_g < -2.0$.

393

394 **Results**

395 **CRISPRi Dataset and Pre-processing**

396 A chemical-genomics dataset was obtained from high-throughput sequencing of a
397 CRISPRi library of *M. tuberculosis* (*Mtb*) that had been treated with several antibiotics. The
398 library consists of 96,700 sgRNAs targeting all 4019 genes in the *Mtb* H37Rv genome [13]. This
399 library was intentionally constructed to focus on probing essential genes (based on prior TnSeq
400 analysis [24]), with a mean of 83 sgRNAs per essential gene, but there are some sgRNAs in each
401 non-essential gene too (mean of 10 sgRNAs per non-essential gene).

402 The library was individually treated with 9 anti-TB drugs (rifampicin, RIF; isoniazid, INH,
403 ethambutol, EMB; vancomycin, VAN; levofloxacin, LEVO; linezolid, LZD; streptomycin, STR;
404 clarithromycin, CLR; bedaquiline, BDQ) to evaluate and validate the CRISPRi system in
405 preparation for target identification for novel inhibitors (from high-throughput screens). These
406 drugs were selected because certain genes are expected to interact for each (based on known
407 mechanisms of action), although additional genes might also exhibit interactions, which could
408 extend our knowledge. We note that some drug targets are members of a complexes; although
409 a drug may bind directly to one subunit, other subunits in those complexes often show similar
410 CRISPRi phenotypes. RIF binds RpoB (RNA polymerase subunit) inhibiting transcription and

411 compensatory mutations are often found in *rpoC* [25], BDQ binds and inhibits AtpE (subunit of
412 the ATP synthase) [26] and *mmpL5* effluxes the drug [27], GyrA and GyrB (subunits of DNA
413 gyrase) would be expected to interact with fluoroquinolones like LEVO [28], EMB targets
414 *embABC* in the arabinogalactan pathway [29, 30], CLR, LZD and STR bind to the ribosome and
415 inhibit translation, which can be protected by rRNA methyltransferases [31-33], VAN binds to
416 peptidoglycan and is expected to interact with genes in the peptidoglycan synthesis pathways
417 [34, 35], and genes such as *inhA*, *katG*, *ahpC*, *ndh*, *mshA* and *cinA* are implicated in the
418 mechanism of action or resistance for isoniazid, an inhibitor of mycolic acid synthesis [36-39].
419 These define selected interactions that would be expected to be observed in a CRISPRi CGI
420 experiment.

421 Samples of the library (pooled cultures) were treated with each of the drugs, with
422 induction of the Sth1 dCAS9 by ATC (anhydrotetracycline), and were sequenced in triplicate at
423 several concentrations for each drug at 2-fold dilutions around the MIC, along with control
424 samples representing the no-drug samples (0 concentration). Three periods of pre-depletion
425 were evaluated: 1, 5, and 10 days (D1, D5, and D10), since it was initially unknown how many
426 days would be optimal for reducing protein expression after induction of CRISPRi. The
427 measurements reported in this experiment are observed counts of sgRNAs, representing the
428 relative proportion of each mutant in the population (pooled culture of CRISPRi mutants).
429 Abundance of a mutant increases or decreases if silencing of the targeted gene causes a change
430 in fitness. Although target proteins are knocked down by inhibiting transcription via CRISPRi,
431 intracellular protein levels are not directly measured in the experiment. Instead, unique
432 nucleotide barcodes representing each sgRNA are amplified from (integrated) plasmids in the

433 cells, sequenced, and counted. The counts reflect the relative abundance of each CRISPRi
434 mutant. Samples were normalized by dividing individual counts for each sgRNA by the sample
435 total (sum over all sgRNAs).

436 In this dataset, prior estimates of sgRNA efficiency were obtained from empirical data by
437 fitting a piecewise-linear equation to fitness over multiple generations, and then using the
438 model for to extrapolate the predicted log-fold change (LFC) each sgRNA at 25 generations [12].
439 The scale for these efficiencies ranged between -25 (highest depletion) and 0 (no depletion). To
440 determine the effect of depletion solely due to the sgRNA (without drug), uninduced samples
441 (in the absence of dCAS9 induction, -ATC) were also sequenced, to provide counts representing
442 mutant abundances in the absence of depletion of targets as an input to the model.

443

444

445 **The CRISPRi-DR model accurately predicts sgRNA abundances from** 446 **sgRNA strength and drug concentration**

447 The CRISPRi-DR model was fitted for all chemical-genetic interaction datasets from Li,
448 Poulton (13) , which included nine drugs tested at three different concentration levels (after 1,
449 5, and 10-days of pre-depletion without drug). The analyses by CRISPRi-DR found a range of
450 tens to hundreds of significant genes for each dataset. Table 1 show a more detailed account of
451 the significant genes founds in these CRISPRi screens by CRISPRi-DR, categorized into depleted
452 (mutant abundance decreases with drug concentration) and enriched (mutant abundance
453 increases with drug concentration).

454

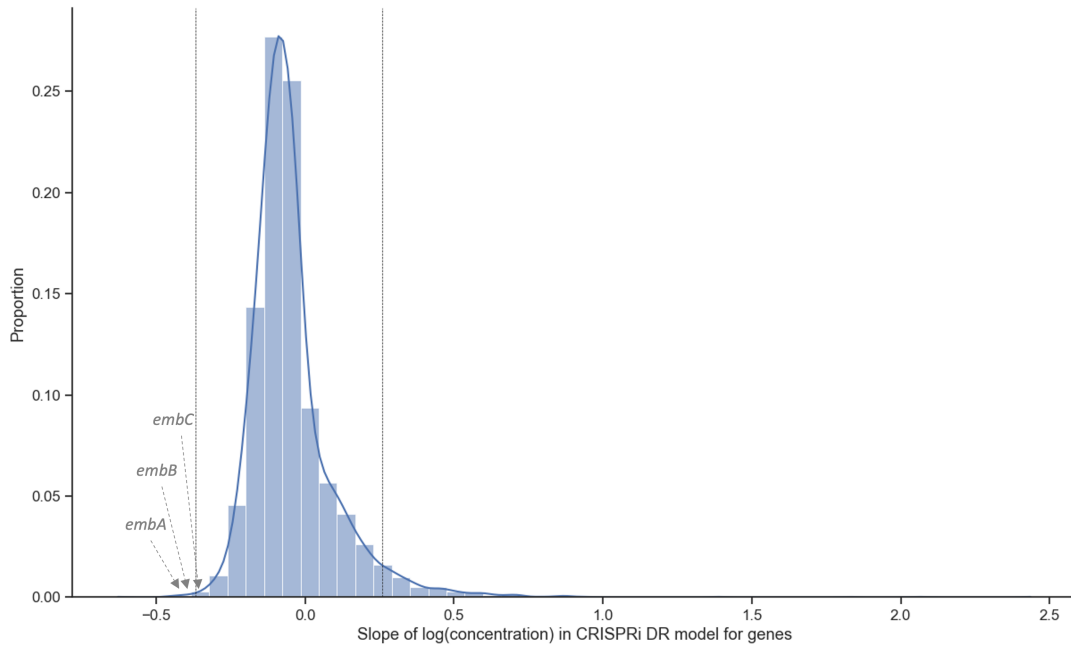
455 **Table 1. Number of Significant Genes found by CRISPRi-DR across the nine drugs CRISPRi**

456 **screen for each of pre-depletion days.**

DRUG	D1		D5		D10	
	Depleted	Enriched	Depleted	Enriched	Depleted	Enriched
BDQ	89	99	121	48	116	37
CLR	182	23	75	89	79	71
EMB	15	160	6	161	51	130
INH	33	57	9	93	16	96
LEVO	80	47	50	50	19	4
LZD	45	123	44	140	54	65
RIF	117	65	165	57	146	53
STR	57	90	44	37	-	-
VAN	193	8	149	26	135	45

457

458 The significant genes identified by CRISPRi-DR generally have coefficients of
459 concentration dependence that are outliers with respect to the rest of the genes. Fig 3 shows
460 the distribution of the slopes calculated for genes in a library treated with EMB (one day of pre-
461 depletion, D1). The threshold for this distribution where $|Z_g| > 2.0$ and adjusted P-value < 0.05 ,
462 is at slope = -0.37 and slope = 0.26 (vertical bars). The 164 total genes in the tails outside the
463 vertical lines are significant genes. These genes include the targets of EMB: *embA*, *embB* and
464 *embC* [29, 30], which have slopes -0.45, -0.43 and -0.32, respectively.



465

466

Fig 3. Coefficients of concentration-dependence from CRISPRi-DR model fitted for EMB

467

D1 (1 day of pre-depletion).

468

The distribution of the slopes of concentration dependence, extracted from the model

469

fit for each gene. The vertical lines are at slope = -0.37 and slope = 0.26. These are the

470

slopes adjusted P-value < 0.05 and the |Z-score| > 2.0. 164 genes have significant slope

471

values, i.e., 164 genes show a significant change in abundance with increasing EMB

472

concentration while accounting for sgRNA strength.

473

474

To evaluate the relative importance of the sgRNA efficiency and drug concentration

475

features to the CRISPRi-DR model, each gene was fit with two ablated models: M_d and M_s . The

476

M_d model contained only log concentration as a predictor: $\log\left(\frac{A_{ijk}}{1-A_{ijk}}\right) = B \cdot \log([D_j]) + C$ and

477

the M_s model only contained sgRNA efficiency as a predictor: $\log\left(\frac{A_{ijk}}{1-A_{ijk}}\right) = B \cdot S_i + C$. In the

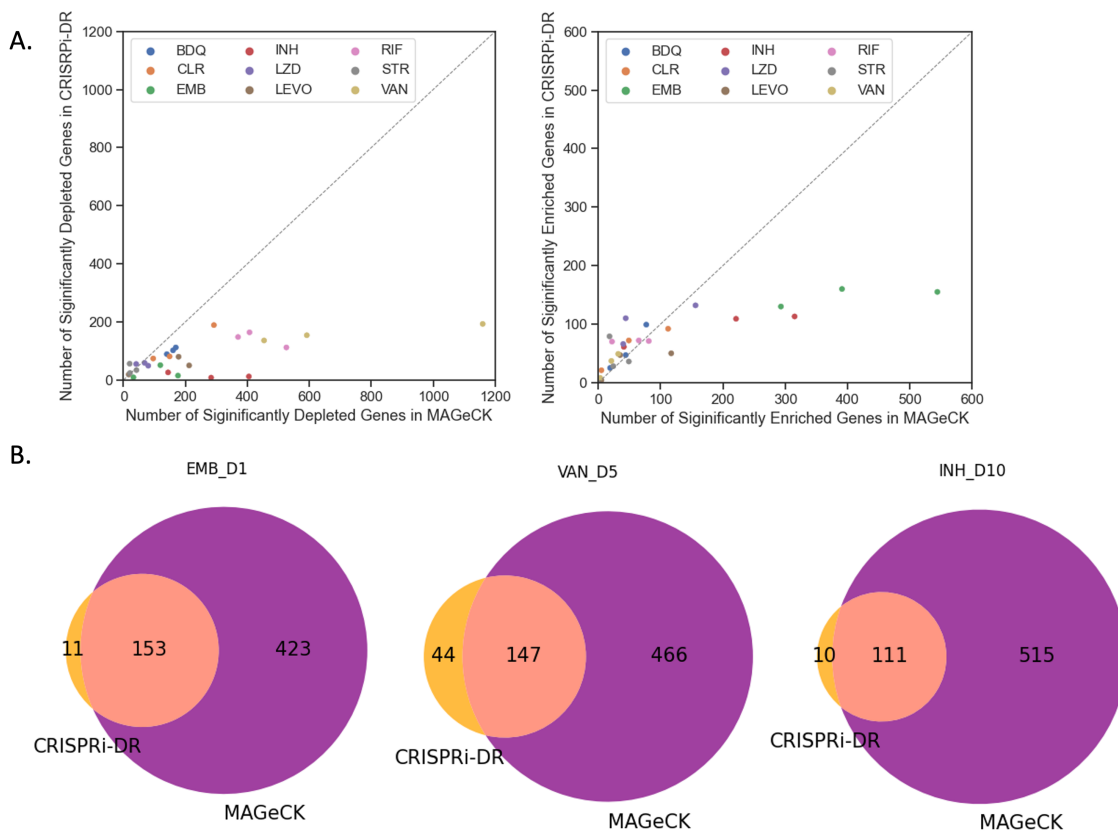
478 EMB D1 experiment, the average r^2 (% variance explained) across all genes in full CRISPRi-DR
479 model is 0.43. Comparatively, the average r^2 is 0.29 for M_s and 0.13 for M_d . *embA* also appears
480 as one of the genes in the M_d set of significant interactors, but the other targets of the drug,
481 *embB* and *embC* do not appear in the sets of significant interactors for either of these ablated
482 models. As a measure of the model quality (goodness of fit), the Akaike Information Criterion
483 (AIC) for the full model in the EMB D1 experiment is 87.6, whereas the AIC of M_d is 300.7 and
484 AIC of M_s is 124.7. The full model has the lowest AIC, indicating it is the best fitting model of
485 the three. The AIC for the model incorporating only drug concentrations but not sgRNA
486 efficiency (M_d) is highest (worst), suggesting that sgRNA efficiency encodes critical information
487 needed for predicting mutant abundance. A Likelihood Ratio Test shows that the differences
488 between these models is highly significant (P-value \ll 0.05; χ^2 distribution using one degree of
489 freedom, since the ablated models each have one parameter less than full model). The r^2
490 values and results of the AIC-based likelihood comparison indicate that sgRNA efficiency
491 contributes strongly to accuracy of the model, and reinforces the importance of including sgRNA
492 efficiency as a term in the CRISPRi-DR model.

493 The improved performance of CRISPRi-DR over the reduced models for EMB extends to
494 the other drugs tested, as seen in Fig. S1. In all the experiments, the number of genes with fits
495 with $r^2 > 0.5$ is the greatest in the full CRISPRi-DR model, and the number of genes with fits that
496 have $r^2 > 0.5$ is greater in model M_s than M_d . This demonstrates that in all conditions, both
497 concentration and sgRNA strength are needed to make accurate estimates of mutant
498 abundance.

499

500 **CRISPRi-DR and MAGeCK have a high concordance of predicted gene-** 501 **drug interactions**

502 Most of the significant CGIs identified by the CRISPRi-DR model were also identified by
503 MAGeCK (MAGeCK-RRA) as reported in Li, Poulton (13), but MAGeCK often identifies many
504 additional genes that are not detected as significant by the CRISPRi-DR model. Although there
505 are some datasets where MAGeCK and CRISPRi-DR detect about the same number of significant
506 interactions, as shown in Fig 4A and the Extended Figure S2 from Li, Poulton (13), there are
507 quite a few datasets where MAGeCK finds substantially more hits than CRISPRi-DR, such as VAN
508 D1, where MAGeCK finds over 1066 significantly depleted genes (even with the filter of $|LFC| > 1$
509 applied), whereas CRISPRi-DR finds only 196 significant interactors. As seen in the Venn
510 diagrams in Fig 4B, there is high overlap of calls made by the two methodologies (enriched and
511 depleted combined). Across all the datasets, an average of 62.2% of genes identified as
512 significant by CRISPRi-DR are also found to be significant by MAGeCK. In the depicted datasets
513 in Panel B, nearly all the calls made by CRISPRi-DR overlap with those made by MAGeCK.
514 However, MAGeCK makes quite a substantial number of calls (significant interacting genes) that
515 are not found by CRISPRi-DR. Additional details of the overlap of significant interacting genes in
516 MAGeCK and CRISPRi-DR can be found in Table S3.



517

518 **Fig 4. Comparison of significant interactions found by CRISPRi-DR and MAGeCK.** (A) The

519 points in the plots are the analyses of CRISPRi screens by both MAGeCK and CRISPRi-DR,

520 colored by drug treatment. The left plot compares the depleted hits called by the two

521 methodologies and the right plot compares the enriched hits called by the two

522 methodologies. The number of hits (both enriched and depleted) are slightly greater in

523 MAGeCK than in the CRISPRi-DR model. (B) Venn Diagram of significant genes, both

524 depleted and enriched, found by CRISPRi-DR and MAGeCK for select drug-treated libraries.

525 The genes identified by CRISPRi-DR are primarily a subset of the hits found by MAGeCK.

526

527 **CRISPRi-DR model correctly detects genes known to interact with anti-**
528 **tubercular drugs.**

529 When genes are ordered by coefficients of the slope representing the dependence of
530 abundance on drug concentration from the CRISPRi-DR model, genes known to affect the
531 potency of the anti-mycobacterial drug tested are ranked highly, as expected (Table 2). The
532 more positive a gene's coefficient is, the higher the gene's enrichment ranking, and the more
533 negative a gene's coefficient is, the higher it's depletion ranking.

534

535 **Table 2: Ranking of Select Genes using the CRISPRi-DR model in 1 Day pre-depletion of treated**
 536 **libraries.**

<i>Drug</i>	<i>Gene</i>	<i>D1 Depletion Ranking</i>	<i>D1 Enrichment Ranking</i>
<i>BDQ</i>	<i>atpA</i>	11	4022
<i>BDQ</i>	<i>atpB</i>	6	4027
<i>BDQ</i>	<i>atpC</i>	51	3982
<i>BDQ</i>	<i>atpD</i>	14	4019
<i>BDQ</i>	<i>atpE</i>	25	4008
<i>BDQ</i>	<i>atpF</i>	9	4024
<i>BDQ</i>	<i>atpG</i>	12	4021
<i>BDQ</i>	<i>atpH</i>	8	4025
<i>BDQ</i>	<i>mmpL5</i>	2	4031
<i>CLR</i>	<i>RVBD3579c</i>	40	3993
<i>CLR</i>	<i>erm(37)</i>	1	4021
<i>EMB</i>	<i>embA</i>	2	4031
<i>EMB</i>	<i>embB</i>	3	4030
<i>EMB</i>	<i>embC</i>	19	4014
<i>INH</i>	<i>inhA</i>	3	4030
<i>INH</i>	<i>ahpC</i>	2	4031
<i>INH</i>	<i>cinA</i>	5	4028
<i>INH</i>	<i>katG</i>	4031	2
<i>INH</i>	<i>ndh</i>	4028	5
<i>INH</i>	<i>mshA</i>	4025	8
<i>LEVO</i>	<i>gyrA</i>	4012	21
<i>LEVO</i>	<i>gyrB</i>	4021	12
<i>LZD</i>	<i>erm(37)</i>	3865	168
<i>LZD</i>	<i>tsnR</i>	4032	1
<i>RIF</i>	<i>rpoB</i>	94	3939
<i>RIF</i>	<i>rpoC</i>	147	3886
<i>STR</i>	<i>RVBD2477c</i>	4021	12
<i>STR</i>	<i>gidB</i>	4022	11

537

538 For each drug, the CRISPRi-DR model is run on each gene (using data from D1). The coefficient
539 for the slope of concentration dependence (β_c) is extracted from the fitted regression and used
540 to rank the genes both in increasing order (for depletion) and inversely (for enrichment). Green
541 reflects results consistent with expectations based on knowledge of known gene-drug
542 interactions

543 Genes that encode the target of a drug would typically be expected to have a high depletion
544 rank, i.e., show a negative slope, indicating that as concentration increases, abundance for the
545 given depletion-mutant decreases. This can be seen in S1 Table, in the ranking for genes using
546 the CRISPRi-DR model. These genes rank the highest in D1 and not as well in D10. This can be
547 attributed to the fact that, after 10 days of pre-depletion, these mutants are already quite
548 depleted, even at concentration 0, increasing noise, and making it difficult to pick up on
549 concentration-dependent signals (further depletion). Therefore, the ranking of relevant genes in
550 D1 was assessed in this analysis (Table 2).

551 For isoniazid (INH), there are multiple relevant genes identified by CRISPRi-DR, including
552 *inhA*, *ahpC*, *ndh* [40], and *katG* [41]. *inhA* (enoyl-ACP reductase) is an essential gene in mycolic
553 acid pathway that is the target of INH, and AhpC (alkyl hydroperoxide reductase) responds to
554 the oxidative effects of isonicotinic radicals in the cells, MshA is a protein involved in synthesis
555 of mycothiol, which helps maintain redox balance [39], and CinA is a NADH metabolizing protein
556 that can hydrolyze the isoniazid-NAD adduct [38]. Therefore, as dosage of the drug increases,
557 the abundances of the mutants of these genes should decrease. These genes are in the top 10
558 highest ranked depletion genes for INH (see Table 2). In contrast, *katG* and *ndh* are found
559 among the top 5 enriched hits, exhibiting increased survival when the proteins are depleted.

560 KatG (catalase) is the activator of INH, and the most common mutations in INH-resistant strains
561 occur in *katG*, decreasing activity [42]. *Ndh* (type II NADH reductase) mutants have also been
562 shown to decrease sensitivity to INH by shifting intracellular NADH levels (needed for INH-NADH
563 adduct formation), and mutations in *ndh* have been shown to be defective in target enzyme
564 (NdhII) activity [40], which is consistent with the observation in the CRISPRi data that depletion
565 of *ndh* leads to increase survival in the presence of INH. Similarly, *mshA* is highly enriched,
566 consistent with mutations found in resistant mutants.

567 For EMB, *embA*, *embB*, and *embC* (subunits of the arabinosyltransferase, target of
568 ethambutol, EMB) rank within the top 100 depleted genes for all three pre-depletion conditions
569 [29, 30]. However, interactions with the other genes in the arabinogalactan pathway, like *ubiA*
570 (which sometimes acquires resistance mutations [43]), were not observed.

571 In RIF, *rpoB* and *rpoC*, subunits of the core RNA polymerase, are ranked within the top 150
572 genes. Significant negative interacting genes for RIF also include many cell wall related genes
573 such as *ponA2*, *rodA*, *ripA*, *aftABCD*, *embABC*, etc., consistent with recent studies that show RIF
574 exposure (or mutations in *rpoB*) leads to various cell wall phenotypes [44-46]. Similarly, the
575 target of bedaquiline (BDQ), the FOF1 ATP synthase (which includes 8 subunits encoded by
576 *atpA-atpH*, of which AtpE is the one bound by BDQ) [26], and *mmpL5*, which can efflux the drug
577 [27], are ranked within the top 40 depleted genes in BDQ.

578 The significantly interacting genes in vancomycin (VAN) involve many genes in the cell
579 wall/membrane/envelope biogenesis pathway (as defined by in COG pathways [47]) (adjusted P-
580 value for pathway enrichment = 0.0004 using Fisher's Exact Test). This follows previous studies

581 that show cell wall genes are targets of vancomycin [48, 49], which binds to peptidoglycan in
582 the cell wall.

583 In levofloxacin (LEVO), CRISPRi mutants of *gyrA* and *gyrB* (subunits of the DNA gyrase, the
584 target of fluoroquinolones) are also observed to be enriched. The reason that depletion of this
585 drug target leads to enrichment of mutants (hence a growth advantage, rather than the
586 expected growth impairment) is likely due to reduced generation of double-stranded breaks in
587 the DNA and other toxic intermediates as a side-effect of inhibiting the gyrase, an effect that has
588 been observed in *E. coli* [50].

589 For clarithromycin (CLR), an inhibitor of translation, *Rv3579c* and *erm(37)* are observed as
590 hits. *Erm(37)* adds a methyl group on the A2058/G2099 nucleotide in the 23S component of the
591 ribosome, the same site in which clarithromycin binds [51]. This natively increases tolerance to
592 CLR in *Mtb*. As this gene is depleted, CLR has greater opportunity to bind, reducing the cells'
593 natural tolerance to the drug. Consistent with this observation, *erm(37)* has a depletion rank of
594 #1 in the CLR D1 condition. *Rv3579c* is another methyltransferase with a similar function that
595 ranks highly (#35) in CLR.

596 In contrast to methylation inhibiting the binding of CLR, there are ribosome
597 methyltransferases in *Mtb*, where methylation facilitates binding of a drug. Mutants for these
598 genes would be expected to show a high enrichment rank in presence of drug. For instance,
599 streptomycin (STR) interferes with ribosomal peptide/protein synthesis by binding near the
600 interaction of the large and small subunits of the ribosome [52]. Two relevant genes that
601 influence the binding of STR include *gidB* and *Rv2477c/ettA*. *GidB* is an rRNA methyltransferase
602 that methylates the ribosome at nucleotide G518 of the 16S rRNA, the position at which STR

603 interacts [33], increasing native affinity for STR. This is consistent with the observation that one
604 of the most common mutations in STR-resistant clinical isolates is loss of function mutations in
605 *gidB* [53]. Rv2477c is a ribosome accessory factor, also known as EttA, which is an ATPase that
606 enhances translation efficiency. It has also recently been shown to bind the ribosome near the
607 P-site (peptidyl transfer center), potentially interfering with binding of aminoglycosides [54],
608 and loss-of-function mutations observed in drug-resistant clinical isolates of *M. tuberculosis*
609 have shown to confer resistance to STR [13]. The ranking of both genes using the CRISPRi-DR
610 model are within the top 12 enriched genes in STR. For linezolid (LZD), relevant genes identified
611 are *erm(37)* and *tsnR*. TsnR is an rRNA methyltransferase, analogous to GidB, and results in
612 tolerance to LZD in a similar manner as GidB does for STR [13]. Following this expectation, *tsnR*
613 has an enrichment ranking of #1 in LZD. Whereas depletion of Erm(37) gives tolerance to CLR, it
614 increases sensitivity to LZD. The nucleotides that Erm(37) methylates in the 23S RNA are
615 proximal in 3D space to where mutations conferring LZD-resistance are found, which both lie in
616 the PTC (peptidyl-transfer center) of the ribosome [55].

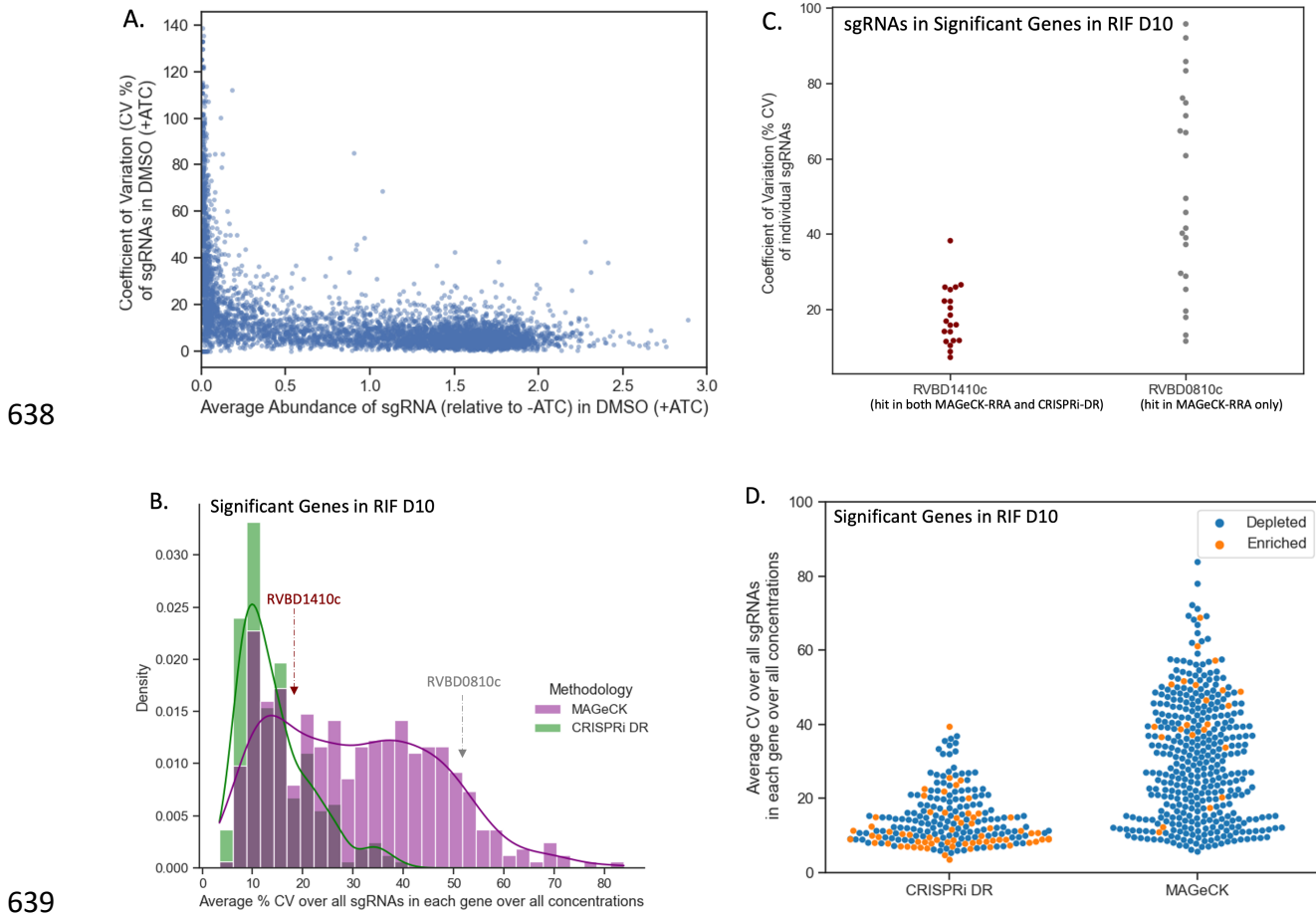
617

618 **The CRISPRi-DR model is less sensitive to noise than MAGeCK**

619 A reason that the CRISPRi-DR model shows lower consistency with MAGeCK (RRA) in some
620 datasets could be due to different sensitivity to noise. There is some noise in these experiments
621 due to variability in sequencing sgRNA counts across multiple concentrations and replicates.
622 This can differentially affect the accuracy of predictions of gene-drug interaction made by these
623 models. Three replicate counts were collected for estimating the relative abundance of each
624 CRISPRi mutant (with a unique sgRNA) in the presence of a drug at a given concentration. The

625 coefficient of variation (CV) can be used to measure the relative consistency of measurements
626 across these observations, which in turn can be used to evaluate the sensitivities of CRISPRi-DR
627 and MAGeCK to noise in the raw data.

628 For each sgRNA s_i the coefficient of variation (CV) was calculated across the relative
629 abundances for the 3 replicates for each concentration (C) in drug (D) ($CV_{D,C,i} = \frac{\sigma(i)}{\mu(i)}$), where
630 $\sigma(i)$ is the standard deviation of the 3 relative abundances in concentration C and $\mu(i)$ is the
631 mean. In Fig 5A, the $CV_{D=DMSO,C=0,i}$ (C of abundances for a random subset of sgRNAs (~5%) in a
632 dCAS9-induced, no-drug condition (concentration 0) is compared to the average abundance. For
633 sgRNAs of medium to high relative abundance (i.e., less depletion), the CV is fairly constant at
634 approximately 10%. However, at low relative (to uninduced) abundances (i.e. higher depletion),
635 CV value increases substantially to over 100%. If a gene contains multiple such sgRNAs with high
636 CV values, then the variation may be misconstrued as a genetic interaction by a methodology
637 that is susceptible to noise.



640 **Fig 5. CRISPRi-DR model shows less sensitivity to noise than MAGECK.** (A) Comparison of
641 average relative abundance and average CV across replicates in no-drug control samples for
642 a sample of sgRNAs: For each sgRNA, we looked at the average CV of sgRNAs in the 3 control
643 replicates against the average abundance of the sgRNA across those replicates. The lower
644 the average abundance, the greater the noise present for the sgRNA. (B) Distribution of
645 average CV of gene for significant genes in MAGECK and significant genes in CRISPRi-DR in
646 RIF D10: The distribution of average CV of significant genes in CRISPRi-DR model is more
647 skewed and has a peak at CV \approx 10%. Although most significant genes in MAGECK show an
648 average CV around 15%, there are quite a few genes with higher average CVs not found
649 significant by the CRISPRi-DR model. (C) Coefficient of Variation (CV) of each sgRNA in two

650 genes with similar number of sgRNAs for a library treated with RIF D10: *Rv1410c* is
651 significant in both methodologies and *Rv0810c* significant in MAGeCK but not in CRISPRi-DR.
652 The majority of CV values for sgRNAs in *Rv1410c* is around 20%. Although both genes have
653 about 20 sgRNAs, *Rv0810c* shows 8 sgRNAs whose CV values exceed 60.5%, which is the
654 maximum CV present in *Rv1410c*. (D) Distribution of average CV for enriched and depleted
655 significant genes in MAGeCK and CRISPRi-DR in a RIF D10 library. This plot shows the
656 distribution plot of Panel B, separated by depletion, and enriched significant genes. The
657 average CV values for significant genes in the CRISPRi-DR model are low for both enriched
658 and depleted genes. As seen in Panel B, significant genes in MAGeCK show low average CV,
659 but they also show high average CV. Although there is a substantially lower number of
660 significantly enriched in MAGeCK, they still show a large amount of noise compared the
661 significantly enriched genes in CRISPRi-DR model.

662
663 The average noise in a gene g for a given drug D can be quantified as the average $CV_{D,C,i}$, for
664 all concentrations C and all sgRNAs in the gene ($\overline{CV}_D(g)$). Therefore, $\overline{CV}_D(g)$ reflects the
665 measure of overall noise present in a gene in a drug D . The distribution of $\overline{CV}_D(g)$ in RIF D10 for
666 the 215 total significant genes (enriched and depleted combined) in the CRISPRi-DR model and
667 in 218 total significant genes (enriched and depleted combined over all concentrations) in
668 MAGeCK can be seen in Fig 5B. The distributions for both methodologies share a peak at about
669 $\overline{CV}_D(g) \approx 10\%$. The distribution of $\overline{CV}_D(g)$ for significant genes in MAGeCK has a fatter tail than
670 the distribution of $\overline{CV}_D(g)$ for significant genes in the CRISPRi-DR model. Fig 5D also shows that
671 the average CV of significant genes found by MAGeCK is much higher than CRISPRi-DR (colored

672 by depleted and enriched) for the RIF D10 screen. In addition to lower CV for significant genes,
673 CRISPRi-DR makes more balanced calls between enriched and depleted, whereas MAGeCK calls
674 are more asymmetric (more depleted than enriched, for this drug). This trend of higher noise in
675 MAGeCK hits is seen not only in RIF D10, but across all the experiments conducted (See S2 Fig).
676 This indicates that although MAGeCK is identifying genes with low noise (like the CRISPRi-DR
677 model), it is also detecting many genes with high noise that the CRISPRi-DR model is not.

678 An example of such a gene is *Rv0810c*. The gene has 22 sgRNAs and has a $\overline{CV}_D(g)$ value
679 (average CV over sgRNAs in a gene) of 51.4%, one of the highest measures in the RIF D10
680 experiment. In RIF D10, it is reported to be significantly depleted only in MAGeCK and not in the
681 CRISPRi-DR model. The dispersion of the CV values of the sgRNAs in *Rv0810c* are compared to
682 those of *Rv1410c* in Fig 5C. *Rv1410c* has 20 sgRNAs, an $\overline{CV}_D(g)$ of 16.3% and is reported to be
683 significantly depleted in both MAGeCK and the CRISPRi-DR model. Although both genes have
684 some sgRNAs with low CVs (below 40%), *Rv0810c* shows 8 sgRNAs with CVs of at least 60.5%,
685 which is the maximum CV of sgRNAs in *Rv1410c*. The CRISPRi-DR model considers the
686 abundances at all concentrations, whereas MAGeCK compares each concentration to the
687 baseline independently. Therefore, if sgRNAs have a high CV value at a particular concentration,
688 they can be picked up as a significant genetic interaction by MAGeCK. The average relative
689 abundance for the 3 replicates at concentration 0 for all sgRNAs in *Rv0810c* is 0.19, whereas the
690 average relative abundance in *Rv1410c* for the same is 1.08. As Fig 5A shows, *Rv0810c* falls in
691 the low abundance/high noise section of the graph, with an average sgRNA no-drug CV of
692 47.9%, whereas *Rv1410c* falls in the low noise section of the graph, with an average sgRNA no-
693 drug CV of 11.2%. This demonstrates that MAGeCK reports genes such as *Rv0810c* with low

694 abundances resulting in a large $\overline{CV}_D(g)$, which the CRISPRi-DR model does not, i.e., MAGeCK is
695 more susceptible to noise than the CRISPRi-DR model.

696

697 **Effects of noise on model performance using simulated CRISPRi data**

698 The sensitivity and accuracy of the CRISPRi-DR model, MAGeCK-RRA and MAGeCK-MLE
699 was assessed under different sources of noise using simulated sgRNA counts sampled from the
700 Negative Binomial distribution [56], with means at different concentrations determined by the
701 dose-response model (Eq (3)). sgRNAs with empirical efficiencies sampled from a uniform
702 distribution from -25 to 0 were used to simulate the combined effects of CRISPRi depletion and
703 exposure to a virtual inhibitor at four concentrations (1 μ M, 2 μ M, 4 μ M, and 8 μ M), with three
704 replicates each. The aim was to determine how noise within and between concentrations
705 affects the performance of each method. Detailed information on the simulation is provided in
706 the Supplementary File S1.

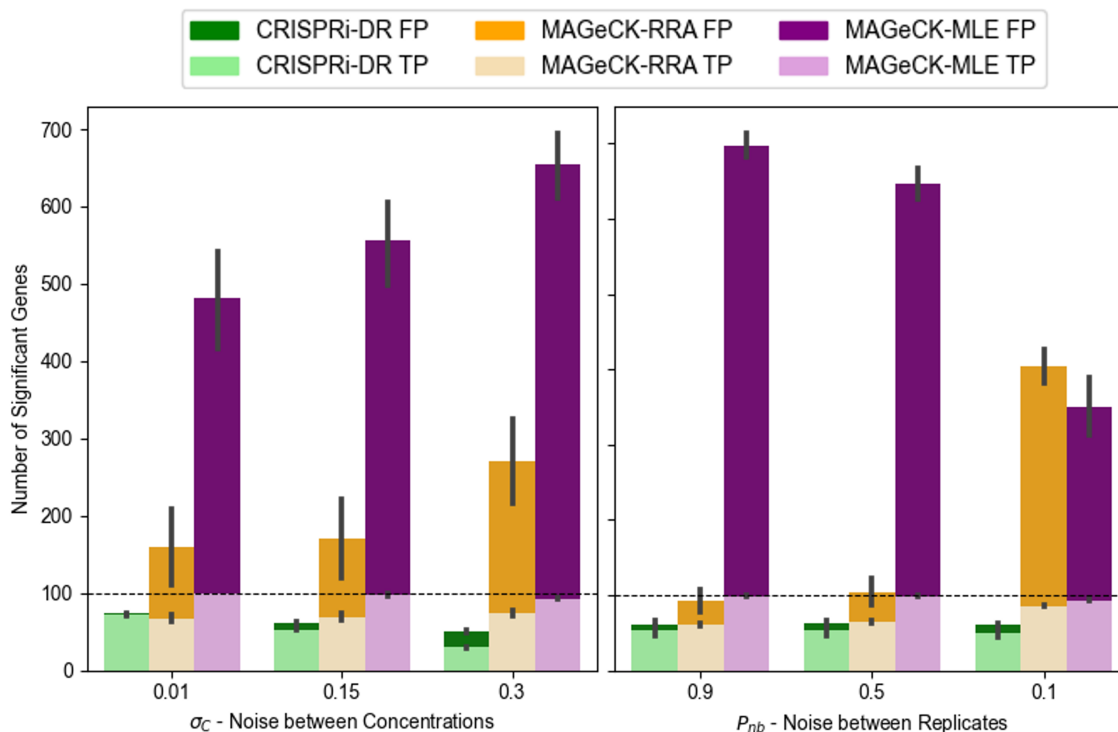
707 Nine datasets (LL, LM, LH, ML, MM, MH, HL, HM and HH) were simulated by varying two
708 noise parameters: variability of abundances *between* concentrations (σ_C), and variability among
709 replicates *within* a concentration (P_{nb} , probability parameter of the Negative Binomial
710 distribution), each with low (L), medium (M), and high (H) setting. A total of 1000 genes was
711 simulated with 20 sgRNAs each. The first 50 genes are chosen as true negative interactions (with
712 a virtual drug), the second 50 as positive interactions, and the last 50 as negative controls (for
713 MAGeCK-RRA and MAGeCK-MLE). For interacting genes, slopes are chosen from a Normal
714 distribution around +0.8 or -0.8, with a standard deviation of 0.2. For non-interacting genes,
715 slopes are chosen from a Normal distribution around 0, with a standard deviation of 0.2.

716 CRISPRi-DR, MAGeCK-RRA and MAGeCL-MLE were run ten times each on these 4 scenarios.
717 MAGeCK was run independently for each drug concentration (2uM, 4uM, 8uM, compared to a
718 no-drug control) and combined using Fisher's method post-hoc, while CRISPRi-DR and MAGeCK-
719 MLE were run on all four concentrations simultaneously.

720 In lowest noise scenario (LL = low noise between concentrations and low noise among
721 replicates), CRISPRi-DR identified 74% of the simulated interacting genes, MAGeCK-RRA
722 identifies 56.5% and MAGeCK-MLE identifies 99.9%. As noise increases, the recall rate of
723 MAGeCK-MLE remains quite high at 88.3% in the highest noise scenario (HH), and MAGeCK-RRA
724 increases to 87.5%. The recall rate of CRISPRi-DR drops down to 30.1%. However, the false
725 positive rate of CRISPRi-DR remains low at 2.2% in this HH scenario, and the false positive rates
726 of MAGeCK-MLE and MAGeCK increase substantially (MLE = 42.5%, RRA = 42.1%), diluting the
727 sets of predicted enriched and depleted genes with non-interacting genes (false positives).
728 Therefore, although CRISPRi-DR identifies less of the true interacting genes in higher noise, it
729 maintains its ability to keep the set of reported interacting genes from being diluted with non-
730 interacting. Across most of the 9 noise scenarios, CRISPRi-DR has higher F1-scores than the
731 other two methods, where $F1\ score = 2 \times \frac{recall \times precision}{recall + precision}$, reflecting a better tradeoff
732 between recall and precision (see Supplemental for more details).

733 The effect of noise on the true and false positive calls made by the methods can be seen
734 in Fig 6, where number of significant genes is plotted for each of the adjusted noise parameters.
735 For MAGeCK-MLE, significant genes were identified as those with adjusted P-value (based on a
736 Wald test) less than 0.05. For MAGeCK-RRA, significant genes were identified as those with
737 adjusted combined P-value less than 0.05 and an |LFC| greater than 1. MAGeCK-RRA is more

738 affected by noise among replicates than between concentrations, as evident by the orange bar
739 for $P_{nb}=0.1$. This is likely a result of stochastic fluctuations of counts at individual drug
740 concentrations that are not necessarily supported at other concentrations. This could help
741 explain the poor performance of MAGeCK-RRA on certain drug-treated screens that may be
742 especially noisy, resulting in many hits, such as in the case of VAN at 1 day pre-depletion; many
743 of these hits could be false positives. Comparatively, CRISPRi-DR and MAGeCK-MLE seem to be
744 more affected by noise between concentrations than noise between replicates, showing lower
745 precision as σ_c increases. Since these methods rely more on increasing or decreasing trends in
746 abundance that must be (at least somewhat) consistent across concentrations, noise between
747 concentrations may make these trends more difficult to identify.



748
749 **Fig 6 Average True Positives (TP) and False Positives (FP) found by CRISPRi-DR,**
750 **MAGeCK-RRA and MAGeCK-MLE as Simulated Noise Increases.** The horizontal dashed

751 line in both panels is the number of total simulated interacting genes (100 total). The
752 parameters in the x-axis are ordered to reflect increasing noise. The leftmost bars of the
753 two plots are the lowest noise and the rightmost bars are the highest noise. MAGeCK-
754 MLE produces a high false positive rate for all scenarios and MAGeCK-RRA is more
755 sensitive to noise among replicates as seen by the orange bar for $P_{nb}=0.1$.

756

757 To assess the impact of performing a CRISPRi screen at multiple drug concentrations on
758 the performance of CRISPRi-DR, MAGeCK and MAGeCK-RRA, we conducted the simulation
759 above with high-noise settings (HH) and varying numbers of drug concentrations (1, 2, or 3) for
760 10 iterations each. The recall of the methods held fairly constant as concentrations were added.
761 However, increasing the number of concentration points caused a significant increase in false
762 positive calls by MAGeCK-RRA from 200 to 400. While MAGeCK-RRA shows susceptibility to
763 false positives when evaluating only a single concentration point, this effect was amplified with
764 more concentrations. This accumulation of errors explains the decrease in precision with
765 additional concentration points. In contrast, CRISPRi-DR is more robust with respect to false-
766 positive errors. By incorporating data from all available concentrations and identifying
767 significant trends, CRISPRi-DR maintains higher precision that does not diminish with the
768 addition of more concentration points. Although MAGeCK-MLE makes many more calls,
769 including false positives, the number of false positives did not increase as concentrations were
770 added, because, like CRISPRi-DR, MAGeCK-MLE incorporates data from all available
771 concentrations.

772

773 **Comparison of CRISPRi-DR to Alternative Methods for CRISPRi**

774 **Analysis**

775 To understand how well CRISPRi-DR performs relative to other CRISPR analysis methods,
776 we applied the following methods on the *M. tuberculosis* CGI data from [13] described above:
777 CGA-LMM [20], MAGeCK-RRA [14], MAGeCK-MLE [15], DrugZ [17], DEBRA [18], and
778 CRISPhieRmix [16]. Each method offers a unique approach to analyzing CRISPRi data. Some of
779 these methods, such as CGA-LMM do not explicitly incorporate multiple sgRNAs per gene or
780 account for differences in sgRNA strength. Other methods, such as DEBRA, MAGeCK-RRA and
781 drugZ, do not explicitly account for different drug concentrations in a CGI experiment, and so
782 they must be run independently on each concentration and the results combined. Only
783 CRISPRi-DR and MAGeCK-MLE incorporate both of these factors in their statistical analysis.

784 The details of applying each method, including parameter settings, handling of negative
785 controls, and merging of results, are described in the Supplement. Several of the methods,
786 including MAGeCK-MLE, produced more significant interactions (in the thousands, in some
787 cases), whereas other methods, like CRISPRi-DR, produced much more focused lists of
788 significant hits for each drug (often less than 100) (see details in the Supplement).

789 To evaluate the accuracy of the predictions by each method, we ranked the genes by
790 significance (usually based on P-value, for most methods) and then generated ROC (Receiver-
791 Operator Characteristic) curves. To define a list of expected hits (i.e. interacting genes) for
792 isoniazid (INH D1, with one day of pre-depletion), we obtained a list of 90 conditionally essential
793 genes from a previously published TnSeq study of *M. tuberculosis* H37Rv exposed to sub-MIC

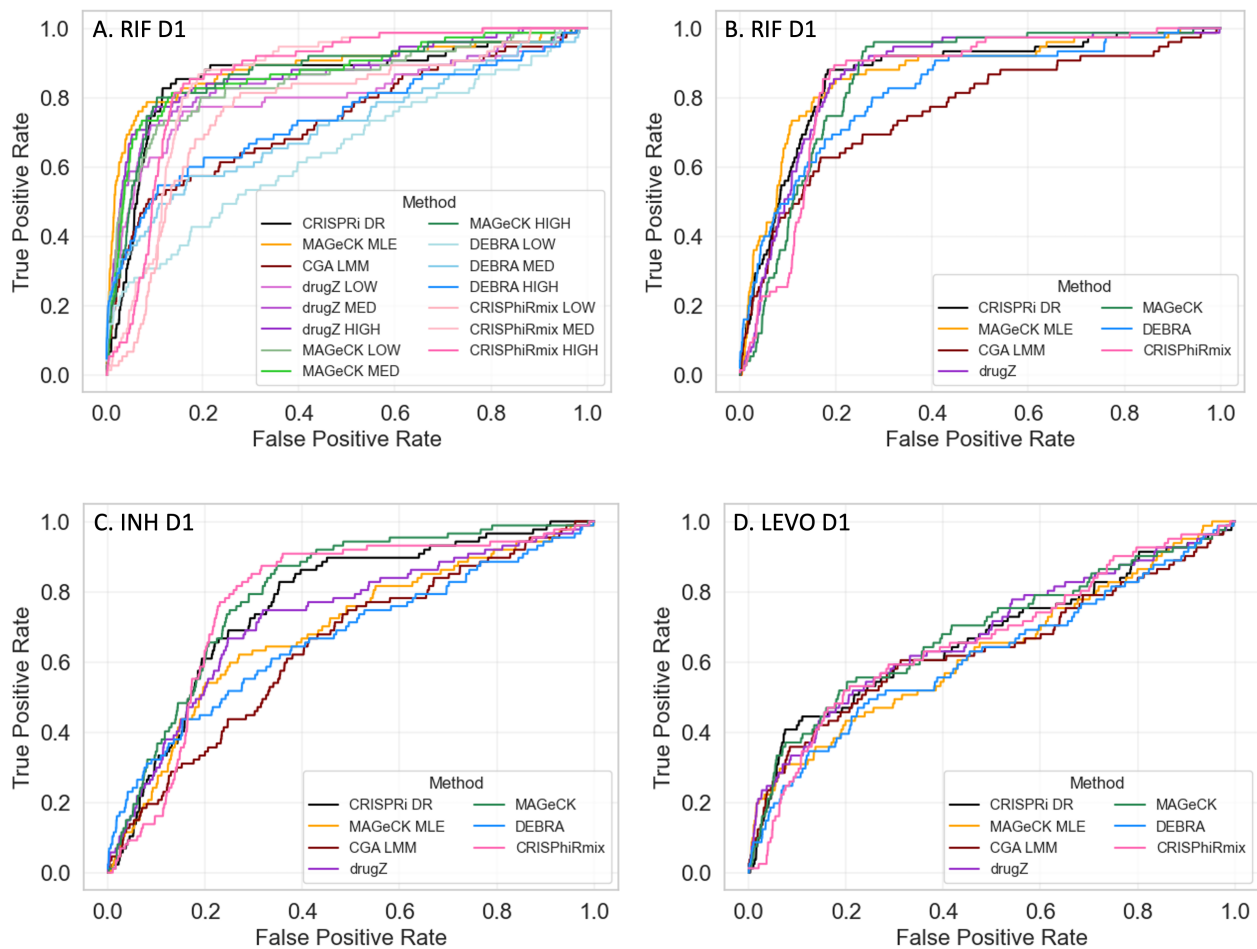
794 concentrations of antibiotics [35]. While changes in essentiality due to knock-out of a gene by
795 transposon insertion are not technically the same as fitness defects resulting from CRISPRi
796 depletion of a target gene, there is substantial overlap between essentiality and vulnerability
797 [12]. Many genes known to play a role in INH resistance (*fabG1*, *katG*, *ndh*, *ahpC*, *cinA*, etc.) are
798 highly interacting (enriched or depleted) in both experiments. Thus, the list of TnSeq
799 conditional essentials serves as a proxy for the genes that are expected to exhibit an interaction
800 effect in the CRISPRi screen (even though, admittedly, not all necessarily will). Importantly,
801 conditional essentiality in this context includes genes whose disruption causes either a growth
802 defect or growth advantage (hypothetically corresponding to depletion or enrichment in a
803 CRISPRi experiment). Similarly, to define a list of expected hits for rifampicin, we used a list of
804 75 conditionally essential genes based on exposure of the TnSeq library to rifampicin, which
805 does not include subunits of the RNA polymerase because they are essential, but includes
806 conditionally essential genes that might play a biological role in tolerating inhibition of
807 transcription [35]. For levofloxacin (LEVO), we used 83 genes in the DNA damage-response
808 pathway (based on the KEGG annotation [57]), plus *rafABC* (recently shown to be involved in
809 DNA damage signaling [58]). Levofloxacin binds to the DNA gyrase (*gyrAB*), which produces a
810 variety of types of damage to DNA, including double-stranded breaks, and requires several DNA
811 replication and repair mechanisms to survive, such as recombination and the SOS response [59,
812 60]. The genes that will exhibit a chemical-genetic interaction with LEVO are likely to overlap
813 substantially with some of the genes in this DNA damage-response pathway.

814 Each of the CRISPR analysis methods was evaluated using these approximate lists of
815 expected hits for each drug. Since some of the methods were not designed to integrate

816 information from multiple concentrations, the methods were initially evaluated by analyzing
817 each concentration (LOW, MED, HIGH) of a given drug independently. Unsurprisingly, the ROC
818 curves showed considerable dispersion of performance (Fig 7A), which was a consequence of
819 both the method and concentration used (expected interactions were often not well-detected
820 at low drug concentrations). Therefore, to make fairer comparisons to methods like CRISPRi-
821 DR, CGA-LMM, and MAGeCK-MLE, we combined the results of each of the other methods over
822 multiple concentrations by using Fisher's method [61] to combine P-values of genes at each
823 concentration (by summing the logs of the P-values, which is similar to taking the geometric
824 mean) and using this to re-rank the genes. This strategy for combining results from multiple
825 concentrations produced more uniform ROC curves for all the methods, as illustrated in Fig 7B.
826 For methods which required a single set of counts per gene, like DEBRA and CGA-LMM, the
827 most efficient sgRNA was chosen per gene.

828 When the results for different concentrations were combined using Fisher's method,
829 many of the methods exhibited reasonably good performance, ranking expected hits highly (Figs
830 8b-d). For example, for INH, 50% of the expected interactions were ranked in roughly the top
831 20% of all genes by most of the methods, and for RIF, the identification of expected interactions
832 (based on TnSeq) was even better (producing higher rankings of expected hits). For LEVO, the
833 ROC curves show lower AUCs for all of the methods, probably due to the fact that not all the
834 genes in the DNA damage response pathway are required to tolerate exposure to
835 fluoroquinolones. Though there were some variations in performance from drug to drug,
836 indicating that differences in performance were drug-specific, the overall performance was
837 matched fairly well, as quantified by the AUC values in Table 3. In particular, the performance of

838 CRISPRi-DR, while not uniformly the best, was comparable to that of the other methods
839 evaluated. It is notable methods such as CGA-LMM and DEBRA that do account for multiple
840 sgRNAs often had the worst performance (lowest AUC values). The similarity in performance
841 suggests that genes that exhibited CGIs (enrichment or depletion, at least at some
842 concentration) in this experiment were easily detected by all the methods evaluated, despite
843 their different analytical frameworks. Although the AUC values for all the methods were
844 comparable, the other methods often reported many more false positives than CRISPRi-DR.
845 CRISPRi-DR tends to have slightly lower recall but much higher precision than the other
846 methods (see Supplemental Table S2), suggesting it makes more conservative calls (see
847 Supplement). However, it has the highest F1-scores in nearly all drug screens evaluated, which
848 reflects the best tradeoff of recall and precision.
849



850

851

852

853

854

855

856

857

858

859

860

Fig 7 ROC Curves for RIF, INH and LEVO with 1 day pre-depletion. Using expected

interactions derived from TnSeq studies [35] (INH and RIF) and the DNA-damage

pathway (for LEVO), ROC Curves are plotted for CRISPRi-DR and 6 other CRISPR analysis

methods. A) For methods that do not take concentration into account (MAGeCK, drugZ,

DEBRA and CRISPhieRmix), each concentration (LOW, MED, HIGH) was analyzed

independently, producing distinct ROC curves. B-D). For methods that do not take

concentration into account, results of the 3 concentrations were combined using Fisher's

method for combining P-values.

861 **Table 3. AUC values for 7 CRISPR analysis methods, showing comparative performance**
 862 **on 3 datasets (drug treatments, with 1 day of pre-depletion), based on the ROC curves**
 863 **in Figure 7.**

864

	INH D1 AUCs	RIF D1 AUCs	LEVO D1 AUCs
Definition of Hits:	90 TnSeq conditional essentials (Xu et al, 2017)	75 TnSeq conditional essentials (Xu et al., 2017)	83 genes in DNA damage response pathway (KEGG)
CRISPRi-DR	0.767	0.850	0.669
CGA-LMM	0.641	0.765	0.638
MAGeCK-RRA	0.799	0.855	0.684
MAGeCK-MLE	0.683	0.865	0.629
drugZ	0.726	0.866	0.678
DEBRA	0.665	0.822	0.615
CRISPhieRmix	0.771	0.844	0.666

865

866

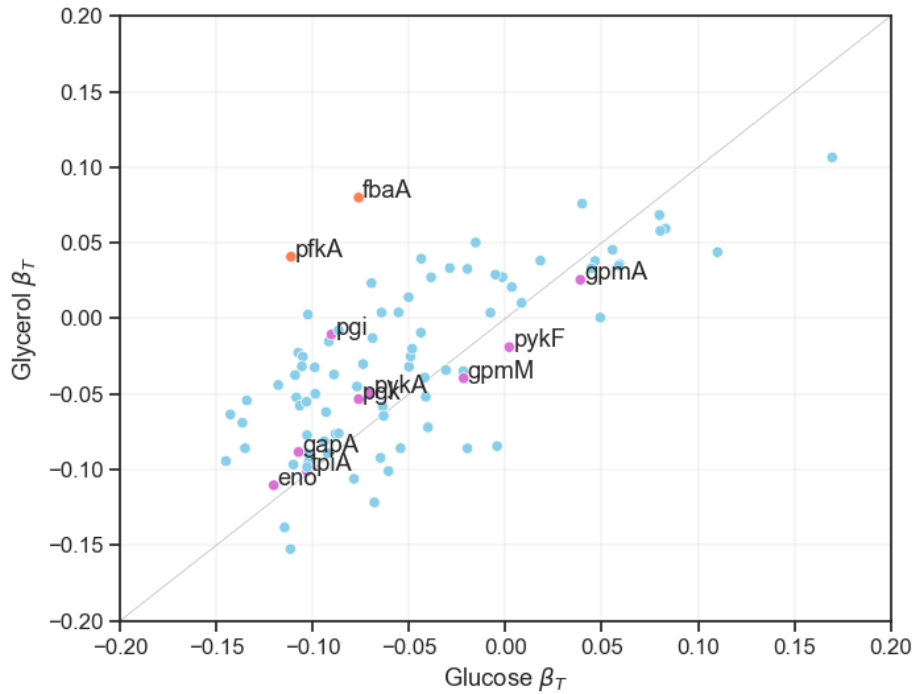
867 **Analysis of CRISPRi Data for *E. coli* Genes Required for Growth on**
 868 **Different Carbon Sources**

869 To illustrate the application of the CRISPRi-DR method to other datasets, we re-analyzed
870 the data from a CRISPRi library in *E. coli* that was used to investigate differential requirements
871 for growth on glycerol versus glucose as a carbon source [11]. While this is not technically a
872 chemical-genetics experiment, the data included multiple time points. The growth curves of
873 CRISPRi knock-down mutants (depletion over time) follows sigmoidal behavior very analogous
874 to dose-response curves for antibiotic exposure (depletion with increasing concentration).
875 Furthermore, while only 88 genes were analyzed instead of a whole-genome screen, this
876 dataset is suitable for analysis by CRISPRi-DR because multiple unique sgRNAs were synthesized
877 for each gene (68 per gene on average), spanning a range of efficiencies (which were quantified
878 by fitting growth data to a logistic curve).

879 We ran CRISPRi-DR on this data for each carbon source independently (fitting the model
880 to 7 timepoints for glucose, 5 for glycerol) (see Supplemental Material for additional details).
881 Many genes exhibited significant depletion effects (reduced fitness), because many of the 88
882 genes were essential for growth (on either carbon source). However, when the coefficients of
883 the time parameter from the CRISPRi-DR analysis were plotted as a scatter plot between the
884 carbon sources, two genes stood out as being preferentially required for growth on glucose
885 (highlighted in orange in Fig 8, most divergent from the diagonal): *fbpA* (fructose biphosphate
886 aldolase) and *pfkA* (phosphofructokinase). These genes are well-known examples required for
887 preliminary steps in glycolysis but not for incorporation of glycerol, and were identified in the
888 analysis by [11]. Additional metabolic genes needed for growth on both carbon sources are
889 observed to lie along the diagonal. This demonstrates that the CRISPRi-DR method can be
890 applied to other datasets, including those not explicitly designed for chemical-genetics. The

891 modified dose-response model nicely incorporates the simultaneous effects of time and the
892 variable efficiency of sgRNAs on mutant abundance.

Coefficients of Time Dependence by CRISPRi-DR in Glucose and Glycerol



893

894 **Fig 8 Coefficients of time dependence from CRISPRi-DR models fit for glucose and**
895 **glycerol *E. coli* datasets.** Each point in the scatterplot represents the coefficients of time
896 dependence of a gene from the fit of the two models (glucose and glycerol). Individually,
897 the gene show a range of growth defect over time, but the coefficients for most genes
898 are equally negative for both conditions, except for a few outliers. The genes colored
899 fuchsia are involved in both gluconeogenesis and glycolysis, hence, as expected, have
900 similar time dependence coefficients in both carbon sources. The points farther away
901 from this line, the orange labeled points (*pfkA* and *fbaA*), are genes involved in glycolysis
902 but not gluconeogenesis and, as expected, they have more negative coefficients in
903 glucose than in glycerol.

904

905 **Discussion**

906 There are a variety of ways to use CRISPRi technology for probing the biological roles of
907 genes by modulating their expression levels in-situ. While early experiments utilized the
908 intrinsic nuclease activity of the CAS9 to knock-out genes entirely [1-3], more recent approaches
909 have enabled partial knock-down of targets, generally using an inactive CAS9 (dCAS9) to bind to
910 target genes and block transcription [5]. One way of controlling the level of depletion is through
911 manipulating the expression of the dCAS9 itself. However, a second approach to creating
912 variability in levels of target depletion is to utilize multiple sgRNAs of different efficiency. The
913 nucleotide sequence of both the PAM and target-specific parts of the guide RNA can impact the
914 hybridization and recruitment of the dCAS9 [9, 10]. This variability can be useful for gauging or
915 titrating phenotypic effects. Rather than all-or-none responses, one can look for genes whose
916 level of depletion correlates with the phenotype of interest.

917 While CRISPRi libraries can be constructed with multiple sgRNAs per target, most CRISPR
918 analytical methods do not explicitly handle such, and those that do (such as MAGeCK-RRA and
919 CRISPhieRmix) are essentially designed to identify significant genes by focusing on a subset of
920 apparently effective sgRNAs (i.e. allowing for ineffective sgRNAs, which are filtered out for each
921 target). However, sgRNA efficiency can be quantified a priori, such as by running a growth
922 experiment to determine the fitness effect of inducing the depletion of the target gene. If this
923 information is available (collected beforehand), then it can be incorporated into the analysis as a
924 “covariate”, to enable comparison of the impact of treatment conditions on the expected

925 magnitude of the phenotypic effect. We note that sgRNA efficiency is different than predicted
926 strength, because it also depends on the vulnerability of the gene. In an essential gene, some
927 sgRNAs might be more efficient than others. In contrast, typically, all the sgRNAs targeting a
928 non-essential will turn out to be non-efficient (i.e. have 0 growth defect, or relative fitness of
929 around 1), at least under control conditions, since the cells are unaffected by depletion of these
930 proteins and continue to grow at the same rate. However, they might cause growth impairment
931 if expressed in certain stress conditions where they might play a role in survival/tolerance. In
932 fact, in chemical-genetic interaction experiments, variable sgRNA efficiency can be further
933 exploited to identify genes whose level of depletion synergizes with increasing drug
934 concentration. We developed the CRISPRi-DR model with this use case in mind, extending the
935 Hill equation, which quantifies dose-response behavior of a growth inhibitor, to incorporate an
936 extra term representing the relative efficiency of each of the sgRNAs targeting a gene. This
937 approach, however, is not limited to CGI experiments. It can be applied to other treatments
938 that induce a sigmoidal response. For example, in re-analysis of data from the Mathis, Otto and
939 Reynolds (11) paper, we showed the same equation could be adapted for modeling the effect of
940 *E. coli* cultures grown on medium with different carbon sources; the time parameter could be
941 substituted for the concentration, since depletion of essential genes caused a gradual killing
942 with an S-curve shape over time.

943 Therefore, the CRISPRi-DR approach we developed has 3 main requirements. First, the
944 CRISPRi library should contain multiple sgRNAs per target gene. Anecdotal evidence suggests
945 that at least 5 sgRNAs per gene are necessary to maintain overall sensitivity for detecting
946 expected interactions and maximizing AUC (based on experiments where we subsampled a

947 limited number of sgRNAs per screen; see Supplement). Fewer sgRNAs per gene reduced the
948 stability of the regression and increased variance of the fitted parameters (specifically the slope
949 of concentration dependence). Second, ideally, sgRNAs of differing strength should be included.
950 Strength can be predicted from sequence features using various types of trained models [9, 12].
951 This covers both essential and non-essential genes. For essential (or vulnerable) genes, sgRNA
952 efficiency correlates with predicted strength, so this is equivalent to choosing sgRNAs with a
953 range of efficiencies (that create varying growth defects). For non-essential genes, one could
954 choose a set of sgRNAs with a range of predicted strengths, even though they might all turn out
955 to be non-efficient experimentally in standard growth conditions. This diversity could be
956 created by selecting sgRNAs that deviate from the optimal PAM sequence [6], choosing
957 hybridizing sequences of different length or GC content [5, 8], or adding random nucleotide
958 substitutions [10]. Third, the actual efficiency of each sgRNA must be empirically quantified a
959 priori, such as by running a growth experiment and comparing growth rates with and without
960 induction of the dCAS9 (hence, with and without depletion of target genes). These quantities
961 become inputs to the model. The CRISPRi-DR method can be applied to any CRISPRi dataset
962 that meets these requirements. The methodology works best when treatment produces a
963 sigmoidal effect on mutant abundances.

964 Doench, Fusi (9) have proposed several systems for design/optimization of CRISPRi
965 libraries. These were more focused on minimizing off-target effects while maximizing
966 sensitivity for detecting of genuine interactions. They do not give a specific recommendation
967 about how many sgRNAs per gene to select. Their library design guidance is to prefer more
968 efficient sgRNAs (e.g. Rule Set 1 selects top 20% of sgRNAs by empirical efficiency and uses

969 these to build a model to predict sgRNA strength; Rule Set 2 extends this with a machine
970 learning model based on additional sequence features to predict sgRNA strength, and prefers
971 sgRNAs with highest score [9]). This contrasts with our approach, where we advocate selecting
972 sgRNAs with a diversity of efficiencies, since we observed that the sgRNAs that exhibited the
973 most synergy with drug treatments were not always the strongest or weakest, but somewhere
974 in the middle of the range.

975 For application to CGI experiments, the availability of CRISPRi data for multiple sgRNAs
976 of varying strengths for each target gene presents new challenges for statistical analysis. In
977 previous work [20], we showed that regressing the relative abundances of mutants in
978 hypomorph libraries over multiple concentrations of a drug (on log-scale) can be used to
979 improve detection of CGIs. This regression approach captured dose-dependent behavior, i.e.
980 genes whose decreased expression caused either suppressed or enhanced fitness that increases
981 in magnitude with drug concentration (i.e. exhibits a trend, which is important for statistical
982 robustness). The CRISPRi-DR method described in this paper extends this previous work by
983 showing how effects of both drug concentration and sgRNA efficiency can be accommodated in
984 the same model. Ideally, interacting genes would be expected to exhibit synergistic behavior
985 with a drug, where depletion of a target protein induces excess depletion (or enrichment) of
986 the mutants grown in the presence of an inhibitor, and this effect is concentration-dependent
987 (exhibits dose-response behavior).

988 In theory, both CRISPRi depletion of essential genes and exposure to antibiotics should
989 impair growth of CRISPRi mutants (at least for depletion of essential genes). One might expect
990 to observe a depletion effect due to either increasing sgRNA efficiency, or drug concentration,

991 each producing regression "slopes" (in log-transformed space), with slopes for sgRNAs targeting
992 non-essential genes being expected to be flat, regardless of predicted sgRNA strength.
993 However, we observed that sgRNA efficiency and concentration effects are not independent -
994 they interact in a non-linear way. sgRNAs that are too weak do not produce enough depletion
995 of a drug target to cause sensitization, and sgRNAs that are too strong deplete a mutant to such
996 low abundances that concentration-dependent effects are difficult to quantify. Often, there is a
997 "sweet spot", or an intermediate sgRNA strength which maximizes the concentration-
998 dependent effect (which could be different for each gene). Our CRISPRi-DR model incorporates
999 both sgRNA efficiency and drug concentration as parameters, and reproduces the non-linear
1000 interaction between them, where the "slopes" for the effect of drug concentration on relative
1001 abundance of mutants can be larger in magnitude for sgRNAs of intermediate strength, while
1002 being flatter (slopes closer to 0) for sgRNAs of high or low strength. MAGeCK-MLE is the only
1003 other analytical method that take sgRNA efficiencies as an input; in that method, the empirical
1004 measures of efficiency are used to initialize the prior probability that each sgRNA is effective
1005 (assuming each gene is represented by a subset of sgRNAs that are effective and others that are
1006 not), which is combined with other conditional probabilities in a Bayesian framework to
1007 determine the posterior probability of interaction for each gene. However, we observed that
1008 MAGeCK-MLE often reports far more significant interactions than CRISPRi-DR or several other
1009 methods and has lower precision.

1010 In this paper, we showed that this non-linear interaction between sgRNA efficiency and
1011 drug concentration can be modeled using an augmented dose-response equation, in which
1012 terms for both effects are included. By fitting the parameters in this equation to CRISPRi data

1013 from a CGI experiment (normalized mutant abundances from sgRNA counts), one can estimate
1014 the degree to which depletion of a given gene sensitizes cells to an inhibitor, and thereby
1015 identify CGIs. While various computational methods exist for fitting non-linear equations, such
1016 as the Levenberg–Marquardt algorithm [62], we chose to linearize the modified Hill equation by
1017 applying a log-sigmoid transform. The transformation enables us to express the equation in a
1018 linear form, where the parameters (IC_{50} , Hill slopes, etc.) appear as coefficients of linear terms
1019 or constants. Consequently, we can use ordinary least-squares regression (OLS) to fit the model
1020 to the CRISPRi dataset.

1021 Sometimes positive and/or negative controls are included in a CRISPRi experiment [8].
1022 While negative controls can be used in methods like MAGeCK-RRA, CRISPRi-DR is not designed
1023 to use controls explicitly in the statistical analysis of CGIs. Hypothetically, negative controls could
1024 be used in the final filtering step to calculate Z-scores for each gene. Instead of basing the Z-
1025 scores on the mean and standard deviation of slope coefficients in the whole set of genes, they
1026 could be based on the distribution of slope coefficients from the negative controls. While we
1027 tested this idea (using 1750 non-targeting sgRNAs included in the *Mtb* CRISPRi dataset as
1028 negative controls), it resulted in many more genes being labeled as interactions (up to half the
1029 genome). It appears that unrelated genes (not involved in the mechanism of action or
1030 resistance to a drug) often have slightly positive or negative random slopes, due to some source
1031 of noise in the experiment that is unaccounted for. Some genes could exhibit weak phenotypic
1032 effects, conferring slight growth defects or advantages under antibiotic stress, even though they
1033 do not play any direct role in the mechanism of action or resistance to the drug. This is the
1034 reason that we advocate identifying genes that are outliers with respect to the rest of the

1035 population of genes, achieved through the filtering step at the end ($|Zscore| > 2$), instead of just
1036 reporting all genes with slope coefficient statistically different from 0.

1037 We compared CRISPRi-DR to several other analytical methods, including MAGeCK-RRA,
1038 MAGeCK-MLE, DEBRA, CRISPhierRMix, CGA-LMM, and drugZ. Some of these methods
1039 incorporate multiple drug concentrations, while other incorporate sgRNA efficiency as an input
1040 to their models. However, only MAGeCK-MLE incorporates both types of input. The
1041 importance of incorporating both inputs in CRISPRi-DR was demonstrated via an experiment
1042 with ablated models; the model fits (AICs) for each gene were significantly worse for models
1043 that regressed abundances against either drug concentration or sgRNA efficiency alone. For
1044 those methods that do not explicitly combine data from multiple drug concentrations and must
1045 be run on each concentration independently, we employed Fisher's method of combining P-
1046 values to create a merged ranking of genes. Using ROC curves to comparing ranking of expected
1047 interactions, CRISPRi-DR performed comparably to the best of these methods, though method
1048 with the highest AUC differed depending on the drug. This evaluation was facilitated by using
1049 lists of conditionally essential genes from TnSeq experiments (exposure to same drugs) to define
1050 an objective list of expected interactions for each drug for making fair comparisons of
1051 performance. However, a major difference observed among the methods was in the number of
1052 significant interactions detected. Methods like CRISPhierRMix, DEBRA, MAGeCK-RRA, and
1053 MAGeCK-MLE produced hundreds to thousands of hits for each drug, whereas CRISPRi-DR
1054 reported a more conservative list of typically less than a hundred interacting genes. It is likely
1055 that many of the interactions detected by the former methods could be false positives. This was
1056 borne out in simulation experiments, where MAGeCK-RRA, and MAGeCK-MLE exhibited

1057 substantially lower precision than CRISPRi-DR. In both the simulated data and real drug screen
1058 datasets, CRISPRi-DR had the highest F1-scores, reflecting the best tradeoff between precision
1059 and recall compared to other methods. Reducing false positives is important because
1060 experimental validation of hits can be expensive, and follow-up is usually only applied to a
1061 handful of top-ranked genes. Furthermore, we used simulated datasets to explore how noise
1062 within or between drug concentrations could affect both the recall and precision of CRISPRi-DR,
1063 MAGECK-RRA, and MAGECK-MLE. Both types of noise increasingly degrade the recall of all
1064 methods, but noise within concentrations (i.e. sgRNA counts among replicates) seemed to cause
1065 the greatest decrease in precision, especially for MAGECK-RRA. The outlier analysis in CRISPRi-
1066 DR (filtering by Z-score in the last step) partially helps to mitigate this, producing a more focused
1067 list of candidate interactions, and hopefully eliminating genes with small random slopes of
1068 concentration dependence that are not genuine interactions (i.e. false positives).

1069

1070 **Data and Code Availability**

1071

1072 A python-based implementation of the CRISPRi-DR method for analyzing CRISPRi data is publicly
1073 available as part of Transit2: <https://transit2.readthedocs.io/en/latest/>

1074

1075 The output files from analyses of the *Mtb* CRISPRi CGI screens from Li, Poulton (13) using
1076 CRISPRi-DR are available for download at: <https://orca1.tamu.edu/CRISPRi-DR/>

1077

1078 **Acknowledgments**

1079 This work was supported by NIH grant P01 AI143575 (TRI, JR, and DS) and by grant INV-
1080 004761 from the Bill and Melinda Gates Foundation (DS and TRI). The funders had no role in
1081 study design, data collection and analysis, decision to publish, or preparation of the manuscript.

References

1082

1083 1. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna JA, Charpentier E. A programmable dual-
1084 RNA-guided DNA endonuclease in adaptive bacterial immunity. *Science*.

1085 2012;337(6096):816-21. Epub 20120628. doi: 10.1126/science.1225829. PubMed PMID:

1086 22745249; PubMed Central PMCID: PMC6286148.

1087 2. Mali P, Yang L, Esvelt KM, Aach J, Guell M, DiCarlo JE, et al. RNA-guided human genome
1088 engineering via Cas9. *Science*. 2013;339(6121):823-6. Epub 20130103. doi:

1089 10.1126/science.1232033. PubMed PMID: 23287722; PubMed Central PMCID:

1090 PMC3712628.

1091 3. Yang H, Wang H, Shivalila CS, Cheng AW, Shi L, Jaenisch R. One-step generation of mice
1092 carrying reporter and conditional alleles by CRISPR/Cas-mediated genome engineering.

1093 *Cell*. 2013;154(6):1370-9. Epub 20130829. doi: 10.1016/j.cell.2013.08.022. PubMed PMID:

1094 23992847; PubMed Central PMCID: PMC3961003.

1095 4. Jensen TI, Mikkelsen NS, Gao Z, Foßelteder J, Pabst G, Axelgaard E, et al. Targeted
1096 regulation of transcription in primary cells using CRISPRa and CRISPRi. *Genome Res*.

1097 2021;31(11):2120-30. Epub 20210818. doi: 10.1101/gr.275607.121. PubMed PMID:

1098 34407984; PubMed Central PMCID: PMC8559706.

1099 5. Qi LS, Larson MH, Gilbert LA, Doudna JA, Weissman JS, Arkin AP, Lim WA. Repurposing
1100 CRISPR as an RNA-guided platform for sequence-specific control of gene expression. *Cell*.

1101 2013;152(5):1173-83. doi: 10.1016/j.cell.2013.02.022. PubMed PMID: 23452860; PubMed

1102 Central PMCID: PMC3664290.

- 1103 6. Rock JM, Hopkins FF, Chavez A, Diallo M, Chase MR, Gerrick ER, et al. Programmable
1104 transcriptional repression in mycobacteria using an orthogonal CRISPR interference
1105 platform. *Nat Microbiol.* 2017;2:16274. Epub 20170206. doi: 10.1038/nmicrobiol.2016.274.
1106 PubMed PMID: 28165460; PubMed Central PMCID: PMC5302332.
- 1107 7. Peters JM, Colavin A, Shi H, Czarny TL, Larson MH, Wong S, et al. A Comprehensive, CRISPR-
1108 based Functional Analysis of Essential Genes in Bacteria. *Cell.* 2016;165(6):1493-506. Epub
1109 20160526. doi: 10.1016/j.cell.2016.05.003. PubMed PMID: 27238023; PubMed Central
1110 PMCID: PMC4894308.
- 1111 8. Gilbert LA, Horlbeck MA, Adamson B, Villalta JE, Chen Y, Whitehead EH, et al. Genome-
1112 Scale CRISPR-Mediated Control of Gene Repression and Activation. *Cell.* 2014;159(3):647-
1113 61. Epub 20141009. doi: 10.1016/j.cell.2014.09.029. PubMed PMID: 25307932; PubMed
1114 Central PMCID: PMC4253859.
- 1115 9. Doench JG, Fusi N, Sullender M, Hegde M, Vaimberg EW, Donovan KF, et al. Optimized
1116 sgRNA design to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nat*
1117 *Biotechnol.* 2016;34(2):184-91. Epub 20160118. doi: 10.1038/nbt.3437. PubMed PMID:
1118 26780180; PubMed Central PMCID: PMC4744125.
- 1119 10. Hawkins JS, Silvis MR, Koo BM, Peters JM, Osadnik H, Jost M, et al. Mismatch-CRISPRi
1120 Reveals the Co-varying Expression-Fitness Relationships of Essential Genes in *Escherichia*
1121 *coli* and *Bacillus subtilis*. *Cell Syst.* 2020;11(5):523-35.e9. Epub 20201019. doi:
1122 10.1016/j.cels.2020.09.009. PubMed PMID: 33080209; PubMed Central PMCID:
1123 PMC7704046.

- 1124 11. Mathis AD, Otto RM, Reynolds KA. A simplified strategy for titrating gene expression reveals
1125 new relationships between genotype, environment, and bacterial growth. *Nucleic Acids*
1126 *Research*. 2020;49(1):e6-e. doi: 10.1093/nar/gkaa1073.
- 1127 12. Bosch B, DeJesus MA, Poulton NC, Zhang W, Engelhart CA, Zaveri A, et al. Genome-wide
1128 gene expression tuning reveals diverse vulnerabilities of *M. tuberculosis*. *Cell*.
1129 2021;184(17):4579-92 e24. Epub 20210722. doi: 10.1016/j.cell.2021.06.033. PubMed
1130 PMID: 34297925; PubMed Central PMCID: PMC8382161.
- 1131 13. Li S, Poulton NC, Chang JS, Azadian ZA, DeJesus MA, Ruecker N, et al. CRISPRi chemical
1132 genetics and comparative genomics identify genes mediating drug potency in
1133 *Mycobacterium tuberculosis*. *Nat Microbiol*. 2022;7(6):766-79. Epub 20220530. doi:
1134 10.1038/s41564-022-01130-y. PubMed PMID: 35637331; PubMed Central PMCID:
1135 PMCPMC9159947.
- 1136 14. Li W, Xu H, Xiao T, Cong L, Love MI, Zhang F, et al. MAGeCK enables robust identification of
1137 essential genes from genome-scale CRISPR/Cas9 knockout screens. *Genome Biol*.
1138 2014;15(12):554. doi: 10.1186/s13059-014-0554-4. PubMed PMID: 25476604; PubMed
1139 Central PMCID: PMC4290824.
- 1140 15. Li W, Koster J, Xu H, Chen CH, Xiao T, Liu JS, et al. Quality control, modeling, and
1141 visualization of CRISPR screens with MAGeCK-VISPR. *Genome Biol*. 2015;16:281. Epub
1142 20151216. doi: 10.1186/s13059-015-0843-6. PubMed PMID: 26673418; PubMed Central
1143 PMCID: PMC4699372.
- 1144 16. Daley TP, Lin Z, Lin X, Liu Y, Wong WH, Qi LS. CRISPhieRmix: a hierarchical mixture model for
1145 CRISPR pooled screens. *Genome Biol*. 2018;19(1):159. Epub 20181008. doi:

- 1146 10.1186/s13059-018-1538-6. PubMed PMID: 30296940; PubMed Central PMCID:
1147 PMCPMC6176515.
- 1148 17. Colic M, Wang G, Zimmermann M, Mascall K, McLaughlin M, Bertolet L, et al. Identifying
1149 chemogenetic interactions from CRISPR screens with drugZ. *Genome Med.* 2019;11(1):52.
1150 Epub 20190822. doi: 10.1186/s13073-019-0665-3. PubMed PMID: 31439014; PubMed
1151 Central PMCID: PMCPMC6706933.
- 1152 18. Akimov Y, Bulanova D, Timonen S, Wennerberg K, Aittokallio T. Improved detection of
1153 differentially represented DNA barcodes for high-throughput clonal phenomics. *Mol Syst*
1154 *Biol.* 2020;16(3):e9195. doi: 10.15252/msb.20199195. PubMed PMID: 32187448; PubMed
1155 Central PMCID: PMCPMC7080434.
- 1156 19. Doench JG, Hartenian E, Graham DB, Tothova Z, Hegde M, Smith I, et al. Rational design of
1157 highly active sgRNAs for CRISPR-Cas9-mediated gene inactivation. *Nat Biotechnol.*
1158 2014;32(12):1262-7. Epub 20140903. doi: 10.1038/nbt.3026. PubMed PMID: 25184501;
1159 PubMed Central PMCID: PMCPMC4262738.
- 1160 20. Dutta E, DeJesus MA, Ruecker N, Zaveri A, Koh EI, Sasseti CM, et al. An improved statistical
1161 method to identify chemical-genetic interactions by exploiting concentration-dependence.
1162 *PLoS One.* 2021;16(10):e0257911. Epub 20211001. doi: 10.1371/journal.pone.0257911.
1163 PubMed PMID: 34597304; PubMed Central PMCID: PMCPMC8486102.
- 1164 21. Wisner MJ, Lenski RE. A Comparison of Methods to Measure Fitness in *Escherichia coli*. *PLoS*
1165 *One.* 2015;10(5):e0126210. Epub 20150511. doi: 10.1371/journal.pone.0126210. PubMed
1166 PMID: 25961572; PubMed Central PMCID: PMCPMC4427439.

- 1167 22. Wald A. The Fitting of Straight Lines if Both Variables are Subject to Error. *The Annals of*
1168 *Mathematical Statistics*. 1940;11(3):284-300.
- 1169 23. Benjamini Y, Krieger AM, Yekutieli D. Adaptive Linear Step-up Procedures That Control the
1170 False Discovery Rate. *Biometrika*. 2006;93(3):491-507.
- 1171 24. DeJesus MA, Gerrick ER, Xu W, Park SW, Long JE, Boutte CC, et al. Comprehensive
1172 Essentiality Analysis of the *Mycobacterium tuberculosis* Genome via Saturating Transposon
1173 Mutagenesis. *mBio*. 2017;8(1). Epub 20170117. doi: 10.1128/mBio.02133-16. PubMed
1174 PMID: 28096490; PubMed Central PMCID: PMC5241402.
- 1175 25. de Vos M, Muller B, Borrell S, Black PA, van Helden PD, Warren RM, et al. Putative
1176 compensatory mutations in the *rpoC* gene of rifampin-resistant *Mycobacterium*
1177 *tuberculosis* are associated with ongoing transmission. *Antimicrob Agents Chemother*.
1178 2013;57(2):827-32. Epub 20121203. doi: 10.1128/AAC.01541-12. PubMed PMID:
1179 23208709; PubMed Central PMCID: PMC3553702.
- 1180 26. Guo H, Courbon GM, Bueler SA, Mai J, Liu J, Rubinstein JL. Structure of mycobacterial ATP
1181 synthase bound to the tuberculosis drug bedaquiline. *Nature*. 2021;589(7840):143-7. Epub
1182 20201209. doi: 10.1038/s41586-020-3004-3. PubMed PMID: 33299175.
- 1183 27. Kaniga K, Lounis N, Zhuo S, Bakare N, Andries K. Impact of Rv0678 mutations on patients
1184 with drug-resistant TB treated with bedaquiline. *Int J Tuberc Lung Dis*. 2022;26(6):571-3.
1185 doi: 10.5588/ijtld.21.0670. PubMed PMID: 35650698; PubMed Central PMCID:
1186 PMC5241402.

- 1187 28. Mayer C, Takiff H. The Molecular Genetics of Fluoroquinolone Resistance in Mycobacterium
1188 tuberculosis. *Microbiol Spectr*. 2014;2(4):MGM2-0009-2013. doi:
1189 10.1128/microbiolspec.MGM2-0009-2013. PubMed PMID: 26104201.
- 1190 29. Cui Z, Li Y, Cheng S, Yang H, Lu J, Hu Z, Ge B. Mutations in the embC-embA intergenic region
1191 contribute to Mycobacterium tuberculosis resistance to ethambutol. *Antimicrob Agents*
1192 *Chemother*. 2014;58(11):6837-43. Epub 20140902. doi: 10.1128/AAC.03285-14. PubMed
1193 PMID: 25182646; PubMed Central PMCID: PMC4249443.
- 1194 30. Zhang L, Zhao Y, Gao Y, Wu L, Gao R, Zhang Q, et al. Structures of cell wall
1195 arabinosyltransferases with the anti-tuberculosis drug ethambutol. *Science*.
1196 2020;368(6496):1211-9. Epub 20200423. doi: 10.1126/science.aba9102. PubMed PMID:
1197 32327601.
- 1198 31. Mougari F, Bouziane F, Crockett F, Nessar R, Chau F, Veziris N, et al. Selection of Resistance
1199 to Clarithromycin in Mycobacterium abscessus Subspecies. *Antimicrob Agents Chemother*.
1200 2017;61(1). Epub 20161227. doi: 10.1128/AAC.00943-16. PubMed PMID: 27799212;
1201 PubMed Central PMCID: PMC5192163.
- 1202 32. Gan WC, Ng HF, Ngeow YF. Mechanisms of Linezolid Resistance in Mycobacteria.
1203 *Pharmaceuticals (Basel)*. 2023;16(6). Epub 20230524. doi: 10.3390/ph16060784. PubMed
1204 PMID: 37375732; PubMed Central PMCID: PMC10303974.
- 1205 33. Wong SY, Lee JS, Kwak HK, Via LE, Boshoff HI, Barry CE, 3rd. Mutations in gidB confer low-
1206 level streptomycin resistance in Mycobacterium tuberculosis. *Antimicrob Agents*
1207 *Chemother*. 2011;55(6):2515-22. Epub 20110328. doi: 10.1128/AAC.01814-10. PubMed
1208 PMID: 21444711; PubMed Central PMCID: PMC3101441.

- 1209 34. Alam MT, Petit RA, 3rd, Crispell EK, Thornton TA, Conneely KN, Jiang Y, et al. Dissecting
1210 vancomycin-intermediate resistance in staphylococcus aureus using genome-wide
1211 association. *Genome Biol Evol.* 2014;6(5):1174-85. Epub 20140430. doi:
1212 10.1093/gbe/evu092. PubMed PMID: 24787619; PubMed Central PMCID:
1213 PMCPMC4040999.
- 1214 35. Xu W, DeJesus MA, Rucker N, Engelhart CA, Wright MG, Healy C, et al. Chemical Genetic
1215 Interaction Profiling Reveals Determinants of Intrinsic Antibiotic Resistance in
1216 *Mycobacterium tuberculosis*. *Antimicrob Agents Chemother.* 2017;61(12). Epub 20171122.
1217 doi: 10.1128/AAC.01334-17. PubMed PMID: 28893793; PubMed Central PMCID:
1218 PMCPMC5700314.
- 1219 36. Palomino JC, Martin A. Drug Resistance Mechanisms in *Mycobacterium tuberculosis*.
1220 *Antibiotics (Basel).* 2014;3(3):317-40. Epub 20140702. doi: 10.3390/antibiotics3030317.
1221 PubMed PMID: 27025748; PubMed Central PMCID: PMCPMC4790366.
- 1222 37. Vilcheze C, Jacobs WR, Jr. The mechanism of isoniazid killing: clarity through the scope of
1223 genetics. *Annu Rev Microbiol.* 2007;61:35-50. doi:
1224 10.1146/annurev.micro.61.111606.122346. PubMed PMID: 18035606.
- 1225 38. Kreutzfeldt KM, Jansen RS, Hartman TE, Gouzy A, Wang R, Krieger IV, et al. CinA mediates
1226 multidrug tolerance in *Mycobacterium tuberculosis*. *Nat Commun.* 2022;13(1):2203. Epub
1227 20220422. doi: 10.1038/s41467-022-29832-1. PubMed PMID: 35459278; PubMed Central
1228 PMCID: PMCPMC9033802.
- 1229 39. Vilcheze C, Av-Gay Y, Barnes SW, Larsen MH, Walker JR, Glynne RJ, Jacobs WR, Jr.
1230 Coresistance to isoniazid and ethionamide maps to mycothiol biosynthetic genes in

- 1231 Mycobacterium bovis. Antimicrob Agents Chemother. 2011;55(9):4422-3. Epub 20110627.
1232 doi: 10.1128/AAC.00564-11. PubMed PMID: 21709101; PubMed Central PMCID:
1233 PMC3165297.
- 1234 40. Vilcheze C, Weisbrod TR, Chen B, Kremer L, Hazbon MH, Wang F, et al. Altered NADH/NAD+
1235 ratio mediates coresistance to isoniazid and ethionamide in mycobacteria. Antimicrob
1236 Agents Chemother. 2005;49(2):708-20. doi: 10.1128/AAC.49.2.708-720.2005. PubMed
1237 PMID: 15673755; PubMed Central PMCID: PMC547332.
- 1238 41. Hazbón MH, Brimacombe M, Bobadilla del Valle M, Cavatore M, Guerrero MI, Varma-Basil
1239 M, et al. Population genetics study of isoniazid resistance mutations and evolution of
1240 multidrug-resistant Mycobacterium tuberculosis. Antimicrob Agents Chemother.
1241 2006;50(8):2640-9. doi: 10.1128/aac.00112-06. PubMed PMID: 16870753; PubMed Central
1242 PMCID: PMC1538650.
- 1243 42. Bollela VR, Namburete EI, Feliciano CS, Macheque D, Harrison LH, Caminero JA. Detection
1244 of katG and inhA mutations to guide isoniazid and ethionamide use for drug-resistant
1245 tuberculosis. Int J Tuberc Lung Dis. 2016;20(8):1099-104. doi: 10.5588/ijtld.15.0864.
1246 PubMed PMID: 27393546; PubMed Central PMCID: PMC5310937.
- 1247 43. Giri A, Gupta S, Safi H, Narang A, Shrivastava K, Kumar Sharma N, et al. Polymorphisms in
1248 Rv3806c (ubiA) and the upstream region of embA in relation to ethambutol resistance in
1249 clinical isolates of Mycobacterium tuberculosis from North India. Tuberculosis (Edinb).
1250 2018;108:41-6. Epub 20171012. doi: 10.1016/j.tube.2017.10.003. PubMed PMID:
1251 29523326.

- 1252 44. McNeil MB, Chettiar S, Awasthi D, Parish T. Cell wall inhibitors increase the accumulation of
1253 rifampicin in *Mycobacterium tuberculosis*. *Access Microbiol.* 2019;1(1):e000006. Epub
1254 20190320. doi: 10.1099/acmi.0.000006. PubMed PMID: 32974492; PubMed Central
1255 PMCID: PMCPMC7470358.
- 1256 45. Patel Y, Soni V, Rhee KY, Helmann JD. Mutations in *rpoB* That Confer Rifampicin Resistance
1257 Can Alter Levels of Peptidoglycan Precursors and Affect β -Lactam Susceptibility. *mBio.*
1258 2023;14(2):e0316822. Epub 20230213. doi: 10.1128/mbio.03168-22. PubMed PMID:
1259 36779708; PubMed Central PMCID: PMCPMC10128067.
- 1260 46. Campodonico VL, Rifat D, Chuang YM, Ioerger TR, Karakousis PC. Altered *Mycobacterium*
1261 *tuberculosis* Cell Wall Metabolism and Physiology Associated With *RpoB* Mutation H526D.
1262 *Front Microbiol.* 2018;9:494. Epub 20180319. doi: 10.3389/fmicb.2018.00494. PubMed
1263 PMID: 29616007; PubMed Central PMCID: PMCPMC5867343.
- 1264 47. Galperin MY, Wolf YI, Makarova KS, Vera Alvarez R, Landsman D, Koonin EV. COG database
1265 update: focus on microbial diversity, model organisms, and widespread pathogens. *Nucleic*
1266 *Acids Res.* 2021;49(D1):D274-D81. doi: 10.1093/nar/gkaa1018. PubMed PMID: 33167031;
1267 PubMed Central PMCID: PMCPMC7778934.
- 1268 48. Provvedi R, Boldrin F, Falciani F, Palu G, Manganelli R. Global transcriptional response to
1269 vancomycin in *Mycobacterium tuberculosis*. *Microbiology (Reading).* 2009;155(Pt 4):1093-
1270 102. doi: 10.1099/mic.0.024802-0. PubMed PMID: 19332811.
- 1271 49. Soetaert K, Rens C, Wang XM, De Bruyn J, Laneelle MA, Laval F, et al. Increased Vancomycin
1272 Susceptibility in *Mycobacteria*: a New Approach To Identify Synergistic Activity against
1273 Multidrug-Resistant *Mycobacteria*. *Antimicrob Agents Chemother.* 2015;59(8):5057-60.

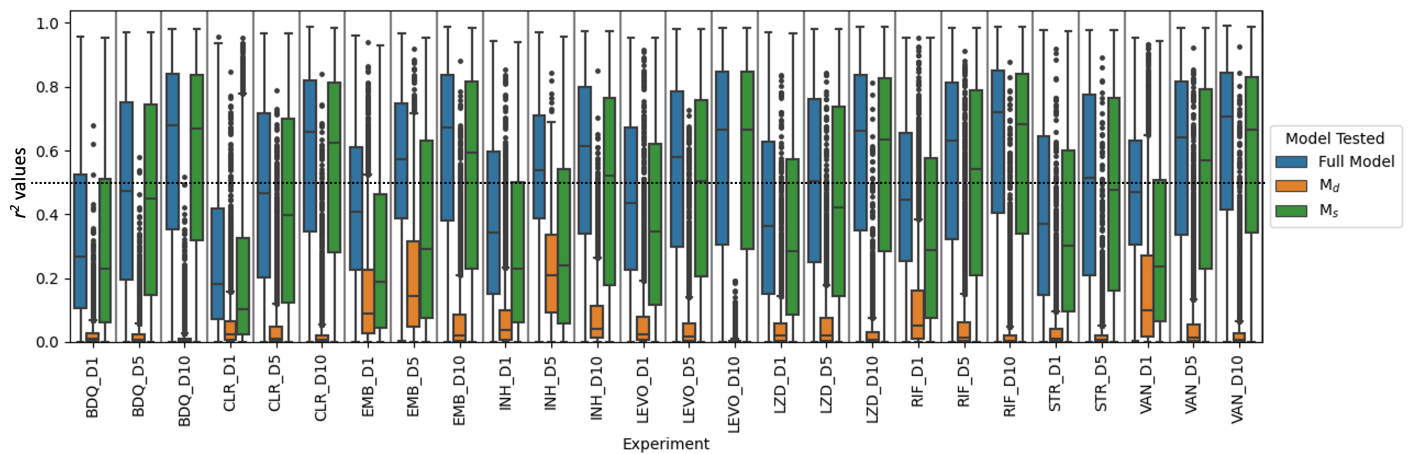
- 1274 Epub 20150601. doi: 10.1128/AAC.04856-14. PubMed PMID: 26033733; PubMed Central
1275 PMCID: PMCPMC4505240.
- 1276 50. Palmer AC, Kishony R. Opposing effects of target overexpression reveal drug mechanisms.
1277 Nat Commun. 2014;5:4296. Epub 20140701. doi: 10.1038/ncomms5296. PubMed PMID:
1278 24980690; PubMed Central PMCID: PMCPMC4408919.
- 1279 51. Hansen JL, Ippolito JA, Ban N, Nissen P, Moore PB, Steitz TA. The structures of four
1280 macrolide antibiotics bound to the large ribosomal subunit. Mol Cell. 2002;10(1):117-28.
1281 doi: 10.1016/s1097-2765(02)00570-1. PubMed PMID: 12150912.
- 1282 52. Chulluncuy R, Espiche C, Nakamoto JA, Fabbretti A, Milón P. Conformational Response of
1283 30S-bound IF3 to A-Site Binders Streptomycin and Kanamycin. Antibiotics (Basel).
1284 2016;5(4). Epub 20161213. doi: 10.3390/antibiotics5040038. PubMed PMID: 27983590;
1285 PubMed Central PMCID: PMCPMC5187519.
- 1286 53. Spies FS, Ribeiro AW, Ramos DF, Ribeiro MO, Martin A, Palomino JC, et al. Streptomycin
1287 resistance and lineage-specific polymorphisms in Mycobacterium tuberculosis gidB gene. J
1288 Clin Microbiol. 2011;49(7):2625-30. Epub 20110518. doi: 10.1128/JCM.00168-11. PubMed
1289 PMID: 21593257; PubMed Central PMCID: PMCPMC3147840.
- 1290 54. Cui ZL, Xiaojun ; Shin, Joonyoung ; Gamper, Howard ; Hou, Ya-Ming ; Sacchettini , James C ;
1291 Zhang, Junjie Interplay between an ATP-binding cassette F protein and the ribosome from
1292 Mycobacterium tuberculosis. Nature Communications. 2022. PubMed Central PMCID:
1293 PMC35064151.
- 1294 55. Madsen CT, Jakobsen L, Buriankova K, Doucet-Populaire F, Pernodet JL, Douthwaite S.
1295 Methyltransferase Erm(37) slips on rRNA to confer atypical resistance in Mycobacterium

- 1296 tuberculosis. *J Biol Chem.* 2005;280(47):38942-7. Epub 20050920. doi:
1297 10.1074/jbc.M505727200. PubMed PMID: 16174779.
- 1298 56. Frazee AC, Jaffe AE, Langmead B, Leek JT. Polyester: simulating RNA-seq datasets with
1299 differential transcript expression. *Bioinformatics.* 2015;31(17):2778-84. Epub 20150428.
1300 doi: 10.1093/bioinformatics/btv272. PubMed PMID: 25926345; PubMed Central PMCID:
1301 PMCPMC4635655.
- 1302 57. Kanehisa M, Furumichi M, Sato Y, Kawashima M, Ishiguro-Watanabe M. KEGG for
1303 taxonomy-based analysis of pathways and genomes. *Nucleic Acids Res.* 2023;51(D1):D587-
1304 D92. doi: 10.1093/nar/gkac963. PubMed PMID: 36300620; PubMed Central PMCID:
1305 PMCPMC9825424.
- 1306 58. Fudrini Olivencia B, Muller AU, Roschitzki B, Burger S, Weber-Ban E, Imkamp F.
1307 *Mycobacterium smegmatis* PafBC is involved in regulation of DNA damage response. *Sci*
1308 *Rep.* 2017;7(1):13987. Epub 20171025. doi: 10.1038/s41598-017-14410-z. PubMed PMID:
1309 29070902; PubMed Central PMCID: PMCPMC5656591.
- 1310 59. Diaz-Diaz S, Recacha E, Machuca J, Garcia-Duque A, Docobo-Perez F, Blazquez J, et al.
1311 Synergistic Quinolone Sensitization by Targeting the recA SOS Response Gene and Oxidative
1312 Stress. *Antimicrob Agents Chemother.* 2021;65(4). Epub 20210318. doi:
1313 10.1128/AAC.02004-20. PubMed PMID: 33526493; PubMed Central PMCID:
1314 PMCPMC8097469.
- 1315 60. Tran T, Ran Q, Ostrer L, Khodursky A. De Novo Characterization of Genes That Contribute to
1316 High-Level Ciprofloxacin Resistance in *Escherichia coli*. *Antimicrob Agents Chemother.*

- 1317 2016;60(10):6353-5. Epub 20160923. doi: 10.1128/aac.00889-16. PubMed PMID:
1318 27431218; PubMed Central PMCID: PMC5038283.
- 1319 61. Mosteller F, Fisher RA. Questions and Answers. The American Statistician. 1948;2(5):30-1.
1320 doi: 10.2307/2681650.
- 1321 62. Helgesson P, Sjostrand H. Fitting a defect non-linear model with or without prior,
1322 distinguishing nuclear reaction products as an example. Rev Sci Instrum.
1323 2017;88(11):115114. doi: 10.1063/1.4993697. PubMed PMID: 29195386.

1324

1325 Supporting Information



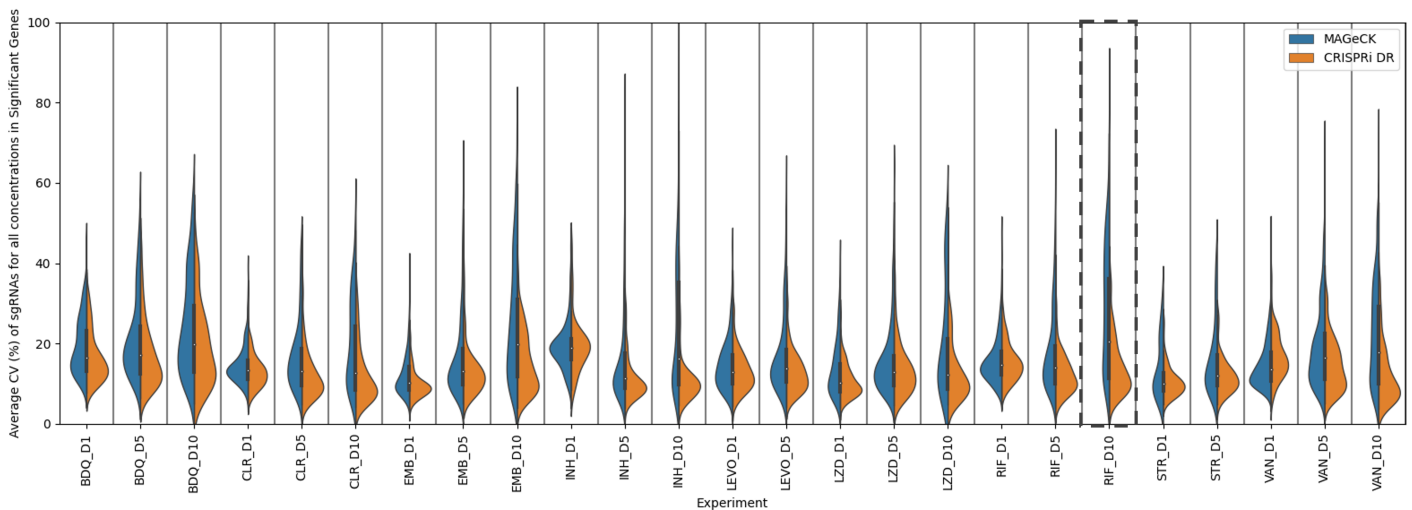
1326

1327 **Fig S1 Evaluation sgRNA strength and log concentration as predictors of CRISPRi-DR**
1328 **through comparison of distribution of r^2 values of full (CRISPRi-DR) and ablated (M_s and M_d)**
1329 **models for each gene in each experiment.**

1330 The horizontal line is where $r^2 = 0.5$. The average r^2 M_s model for all genes across all the
1331 experiments is 0.42, the average r^2 for the M_d model is 0.07. This alongside the Log-likelihood
1332 tests indicate sgRNA strength is the more significant predictor. However, the full CRISPRi-DR

1333 model outperforms both M_d and M_s (average r^2 is 0.50) indicating the inclusion of both sgRNA
1334 strength and log concentration is needed for accurate assessment of significant sgRNA depletion
1335 in a gene in a condition.

1336



1337

1338 **Fig S2 Distribution of average CV of sgRNAs in significant genes (depleted and enriched) in the**
1339 **CRISPRi-DR model and MAGECK.**

1340 In this Fig, we see all the noise distributions for hits in MAGECK and the CRISPRi-DR model for
1341 all experiments. The dashed panel is that of RIF D10. The same distribution of noise of hits can
1342 be seen in Fig 5. The trend seen with RIF D10 is present with all the experiments except LEVO
1343 D10. We see that the CRISPRi-DR model is unimodal with a low CV as the mode, whereas
1344 MAGECK shows significant genes with low average CV values but also a significant amount of
1345 genes with high average CV values. LEVO D10 was left out of this plot due to the low number of
1346 hits in either model.

1347

1348 **Table S1. Ranking of Select Genes using the CRISPRi-DR model in 1 Day, 5 day and 10 Day pre-**
1349 **depletion of treated libraries.**

1350 An extended version of Table 2, where the CRISPRi-DR model is run on each gene for each drug
1351 and pre-depletion day. The coefficient for the slope of concentration dependence (β_c) is
1352 extracted from the fitted regressions and used to rank the genes in both increasing order (for
1353 depletion) and inversely (for enrichment). Green reflects results consistent with expectations
1354 based on knowledge of known gene-drug interactions.

1355

1356 **Table S2. Comparison of significant interactions Identified by CRISPR analysis methods of**
1357 **EMB, INH, LEVO, VAN and RIF CRISPRi screens**

1358 For each drug and pre-depletion day of the selected datasets, all 7 CRISPR methods were run.
1359 For methods that do not account for multiple concentrations, they were run separately for each
1360 concentration and the overall significant interactions are also addressed post-combination of
1361 the individual runs using Fisher's method. The comparison of the significant interactions
1362 identified by the models was evaluated using an objectively defined list of true positives. The
1363 genes identified by Xu, DeJesus (35) were used as the "ground truth" against which the other
1364 model's results were compared. For LEVO, genes in the DNA Damaging pathway are used.
1365 Recall, Precision and F1-score columns are colored such that higher values are more green.

1366

1367 **Table S3. Matrices for comparison of significant interactions Identified by CRISPRi-DR and**
1368 **MAGeCK for each drug and pre-depletion day.**

1369 The table presents the results of CRISPRi-DR and MAGeCK analyses for different drugs and pre-
1370 depletion days. Significant interactions are compared in matrix form. Cells with red font indicate
1371 low overlaps between the interactions found by the two models, while cells with green font
1372 represent high overlaps.

1373

1374 **Supplemental File S1**

1375 We expand on the following four topics from the main text in this document: 1) An assessment
1376 of CRISPRi-DR, MAGeCK and MAGeCK-MLE on datasets with simulated noise, 2) Comparison of
1377 CRISPRi-DR to other analysis methods using CGI datasets, 3) Analysis of *E. coli* CRISPRi screens
1378 using CRISPRi-DR and, 4) The minimum number of sgRNAs recommended per gene in CRISPRi-
1379 DR.

1380