

# shinyDepMap, a tool to identify targetable cancer genes and their functional connections from Cancer Dependency Map data

Kenichi Shimada\*, John A Bachman, Jeremy L Muhlich, and Timothy J Mitchison

Laboratory of Systems Pharmacology and Department of Systems Biology, Harvard Medical School, Boston, MA 02115, USA. \*Corresponding Author and Lead Contact.

Contact Info: kenichi\_shimada@hms.harvard.edu

## Abstract

Individual cancers rely on distinct essential genes for their survival. The Cancer Dependency Map (DepMap) is an ongoing project to uncover these gene dependencies in hundreds of cancer cell lines. To make this drug discovery resource more accessible to the scientific community we built an easy-to-use browser, shinyDepMap (<https://labsyspharm.shinyapps.io/depmap>). shinyDepMap combines CRISPR and shRNA data to determine, for each gene, the growth reduction caused by knockout/knockdown and the selectivity of this effect across cell lines. The tool also clusters genes with similar dependencies, revealing functional relationships. shinyDepMap can be used to 1) predict the efficacy and selectivity of drugs targeting particular genes; 2) identify maximally sensitive cell lines for testing a drug; 3) target hop, i.e., navigate from an undruggable protein with the desired selectivity profile, such as an activated oncogene, to more druggable targets with a similar profile; and 4) identify novel pathways driving cancer cell growth and survival.

## 24 **Introduction**

25 Cancer is a disease of the genome. Hundreds, if not thousands, of driver mutations  
26 cause cancer in different patients (Bailey et al., 2018) and extensive collaborative efforts such  
27 as the Cancer Genome Atlas Program (TCGA) have helped discover them (The Cancer  
28 Genome Atlas Research Network, 2019). Targeted therapies, a type of precision medicine,  
29 aim to treat cancer by selectively killing cancer cells with a specific genotype and spectrum of  
30 driver mutations (Friedman et al., 2015). The underlying hypothesis is that cancers depend on  
31 essential genes that are not the same for all tissues, and that these conditionally essential  
32 genes constitute a druggable dependency—an “Achilles’ heel”—that can be exploited to  
33 develop targeted drugs with minimal toxicity. To achieve this goal, it is important to identify  
34 conditionally essential genes for all cancers. It is also important to group these conditionally  
35 essential genes into functionally related sets to maximize the chance of finding a druggable  
36 target within each set, such as a kinase or other enzyme.

37 The concept of essential vs non-essential genes arose largely from genetic research in  
38 model organisms. Traditionally, it was considered a binary distinction that held across any  
39 genotype. However, loss of a given gene can decrease cell growth without killing the cell, so it  
40 is more realistic to assign a numerical value to the degree of essentiality, *i.e.*, the extent to  
41 which loss of a gene, or inhibition of its product, influences fitness. In cancer, this value may  
42 depend on the genotype, transcriptome, and lineage of the cell. In principle, genes that are  
43 only essential in a few cell types might make better drug targets since inhibiting their function is  
44 less likely to cause toxicity in non-cancer tissues. For example, the epidermal growth factor  
45 receptor is strongly required in certain cancer cells, but not in normal bone marrow stem cells,  
46 making it a potentially good target (Wang et al., 2006).

47       The Cancer Dependency Map (DepMap) is an ongoing project to identify essential  
48 genes across hundreds of cancer cell lines using genome-wide CRISPR and shRNA screens  
49 (Tsherniak et al., 2017; Behan et al., 2019). It has already been used successfully to discover  
50 cancer cells' genetic vulnerabilities (Sandoval et al., 2018; X. Wang et al., 2019). These data  
51 represent a gold mine of useful information for biologists and drug developers, but can be  
52 challenging for non-bioinformaticians to manipulate and interpret.

53       The DepMap portal website (<https://depmap.org/portal>) provides a range of information  
54 for each gene, including 1) the cell lines and lineages dependent on the gene, 2) co-dependent  
55 genes (*i.e.*, other genes whose effects on growth are strongly positively or negatively  
56 correlated with the gene), and 3) basal transcript abundances, copy numbers, and mutations  
57 for the gene. However, the DepMap portal has no native tools to integrate CRISPR and shRNA  
58 or to examine functional relationships among essential genes beyond pairwise comparisons.

59       Here we describe shinyDepMap, a web tool to enable researchers to rapidly determine  
60 the essentiality and selectivity of a given gene across cell lines and to find groups of  
61 functionally related genes with similar essentiality profiles. shinyDepMap integrates data from  
62 both CRISPR and shRNA screens, yielding robust measures of the effects of gene loss on cell  
63 viability. From these combined effect scores we derive two measures for each gene: the  
64 degree to which loss of the gene reduces cell growth in sensitive lines ("efficacy"), and the  
65 degree to which its essentiality varies across lines ("selectivity"). To help researchers identify  
66 potential therapeutic targets we clustered genes with strong efficacy scores into functional  
67 units, many of which represent complexes or biological pathways, as previously reported (Pan  
68 et al., 2018). The results of this analysis are accessible via a simple interactive web-tool at  
69 <https://labsyspharm.shinyapps.io/depmap>.

70

## 71 **Results**

### 72 **Assessment of consistency between CRISPR and shRNA dependency scores**

73       The DepMap project (<https://depmap.org/>) provides two separate pre-processed  
74 genome-wide genetic perturbation datasets for hundreds of cell lines using either shRNA or  
75 CRISPR (Meyers et al., 2017; McFarland et al., 2018). In both datasets, the preprocessed  
76 scores represent the growth effects of knocking the gene down or out, with a strongly negative  
77 value in a particular cell line indicating essentiality. Though the pre-processing algorithms for  
78 the shRNA data take “off-target” genes into account when generating essentiality scores, we  
79 nevertheless expected that the essentiality profiles would differ somewhat between shRNA  
80 and CRISPR due to their distinct mechanisms of reducing gene expression.

81       To assess the consistency between CRISPR and shRNA dependency scores, we first  
82 compared the gene/cell line combinations tested with both methods (15,847 genes in 423 cell  
83 lines, Figure 1 – source data 1) and computed essentiality thresholds for each distribution such  
84 that a dependency score more negative than the threshold is considered essential (Figure 1A;  
85 Methods). These thresholds define the subsets of cell line/gene combinations that are  
86 determined to be essential by either CRISPR or shRNA but not both (areas A and B in Figure  
87 1A). The two methods were somewhat consistent on average, with both methods yielding  
88 approximately normal dependency score distributions with mean zero and a left-skewing tail  
89 corresponding to the subset of essential genes (Pearson correlation: 0.456, Spearman  
90 correlation: 0.201).

91       Despite this concordance, the comparison highlighted differences between the two  
92 methods at the individual gene level. First, CRISPR tends to detect weak to moderate gene



93 deletion effects more sensitively, as evidenced by the greater density of CRISPR-essential  
94 genes above the diagonal in the joint distribution plot (Figure 1A). For example, while both  
95 methods identify RAN, CRISPR identifies CCND1 as more essential (Figure 1B). Second,  
96 some genes were shown essential only by CRISPR or shRNA, but not by the other method  
97 (e.g., FOXD4 and EIF5B; Figure 1B).

98 To better understand these inconsistent dependencies, we used Fisher's exact test to  
99 determine which genes, across their perturbations in all 423 cell lines, were enriched for  
100 inconsistent dependencies. We found that 958 and 20 genes were claimed commonly  
101 essential only by CRISPR or shRNA, respectively (Figure 1C). Notably, these two sets of  
102 inconsistently essential genes are enriched for involvement in distinct pathways: for example,  
103 tRNA metabolic process and mitochondrial translation are overrepresented in the CRISPR-  
104 only set, whereas cytosolic translation initiation is overrepresented in the shRNA-only set  
105 (Figure 1D). This suggests that CRISPR and shRNA have distinct biases in assessing some  
106 genes' essentiality, affecting different classes of genes. While CRISPR is considered to be less  
107 susceptible to off-target effects (Smith et al., 2017) and thus now generally preferred over  
108 shRNA, the results and their therapeutic relevance may also depend on the genes of interest.  
109 For example, EIF5B, a gene involved in translation initiation, is highly conserved throughout  
110 Bacteria, Archaea, and Eukarya, suggesting that it may be essential for most human cells  
111 (Sørensen et al., 2001). However, only shRNA, but not CRISPR, highlighted it as an essential  
112 gene. A method that can combine the two dependency scores would compensate for each  
113 other's artifacts and give more robust scores.

114  
115 **A new dependency score combining CRISPR and shRNA**

116 To summarize both CRISPR and shRNA dependency scores,  $S^C$  and  $S^R$ , we developed  
117 a new dependency score by combining them. Recent studies have shown that similar  
118 approaches are practical (Gilvary et al., 2019; W. Wang et al., 2019). Our new dependency  
119 score,  $S^\theta$ , was computed as the weighted average of the two such that  $S^\theta = \theta S^C + (1 - \theta)S^R$ ,  
120 where  $\theta$  is the mixing ratio of the two scores. It is appropriate to have any  $\theta \in [0,1]$ , but we  
121 selected six values of  $\theta: \theta \in \{0, 0.2, 0.4, 0.6, 0.8, 1\}$  in this study, which are equivalent to mixing  
122 CRISPR and shRNA scores at 0:100 ( $=S^R$ ), 20:80, 40:60, 60:40, 80:20, and 100:0 ( $=S^C$ ). Using  
123 the equation above, we computed  $S^\theta$  for 15,847 genes in 423 cell lines for each  $\theta$  (Figure 2 –  
124 source data 2). The distribution of  $S^\theta$  was located between  $S^C (= S^1)$  and  $S^R (= S^0)$  (Figure 2B).  
125

#### 126 **Efficacy: gene essentiality in a sensitive cell line**

127 We used the combined CRISPR-shRNA gene dependency scores to identify genes that  
128 are either commonly or selectively essential in the cell line panel. This distinction is important  
129 for identifying therapeutic targets because inhibition of commonly essential genes may be toxic  
130 to both cancer cells and normal cells, whereas genes that are selectively essential to particular  
131 cancers may allow for a greater therapeutic window. To capture the therapeutic potential of  
132 selectively essential genes we characterized gene dependency effects with two parameters:  
133 the *efficacy*, which defines the strength of the effect in a sensitive cell line, and the *selectivity*,  
134 which describes the variation of the effect across cell lines.

135 We defined the efficacy  $\mathcal{E}_{G,X}^\theta$  as the  $X$ -th percentile of the distribution of combined  
136 dependency scores  $S^\theta$  for gene  $G$  across all cell lines, denoted  $S_G^\theta$ . For a given  $X$ , we defined  
137 a gene as essential when  $\mathcal{E}_{G,X}^\theta$  is lower than the essentiality threshold  $T_\theta$ , which is determined  
138 from the distribution of  $S^\theta$  for all genes in all cell lines (Figure 2C, top panel). Smaller values of

139  $X$  (lower percentiles) lead to more extreme efficacy values  $\mathcal{E}_{G,X}^{\theta}$  and identify more essential  
140 genes in smaller subsets of cell lines (Figure 2C, bottom panel). We selected  $X=1$  for most of  
141 our analysis. This is equivalent to claiming a gene essential when roughly 5 out of 423 cell  
142 lines show dependence on the gene ( $p=7.8e-5$ , binomial test). We discovered from 3,621 to  
143 5,094 commonly and selectively essential genes from  $S^{\theta}$  with different  $\theta$  (Figure 2D).  
144 Reflecting the inconsistencies between CRISPR and shRNA, only 56% (2,037 genes) of the  
145 essential genes overlapped between  $S^C$  and  $S^R$  (Figure 2D). As for the 958 and 20 genes  
146 claimed commonly essential only by CRISPR or shRNA (Figure 1C), the essential genes  
147 discovered with the same method naturally included all of them; however, the essential genes  
148 discovered with the other method included only 40 to 60% of them, highlighting the  
149 inconsistencies between them (Figure 2E). On the other hand, the combined dependency  
150 scores  $S^{\theta}$ , particularly when  $\theta = 0.2, 0.4$ , and  $0.6$ , provide a more sensitive measure,  
151 discovering most of the essential genes claimed by either method. The first principal  
152 component line between  $S^C$  and  $S^R$  was parallel to the line with  $\theta = 0.66$  in Figure 2A, which  
153 maximizes the variance of  $(S^C, S^R)$ . Among the six lines,  $\theta = 0.6$  (CRISPR: shRNA=60:40), is  
154 most similar to this principal component line. Therefore, we chose  $\theta = 0.6$  or the corresponding  
155  $S^{\theta}$  primarily for the rest of the analysis and compare the performance of different  $\theta$  later.

156

### 157 **Selectivity: the difference in gene essentiality among cell lines**

158 We next defined selectivity, a measure of the cell line dependence of the response to  
159 the loss of a gene. Selectivity implies that gene loss has a widely varying effect across the  
160 population of cell lines, such that the dispersion of the score distribution for a selectively  
161 essential gene would be greater than that for a commonly essential gene. We defined the

162 dispersion of gene  $G$ ,  $\mathcal{D}_{G,X}^\theta$ , as the difference between the  $X$ -th and  $(100-X)$ -th percentiles of  
 163  $S_G^\theta$ , or  $\mathcal{D}_{G,X}^\theta = \mathcal{E}_{G,100-X}^\theta - \mathcal{E}_{G,X}^\theta$ . We found that  $\mathcal{E}_{G,X}^\theta$  and  $\mathcal{E}_{G,100-X}^\theta$  were related linearly for the  
 164 majority of genes, corresponding to non-essential and commonly essential genes, while some  
 165 genes had large positive residuals, corresponding to selectively essential genes (e.g., green vs  
 166 orange in Figure 3A-B). We therefore defined the selectivity  $\mathcal{S}_{G,X}^\theta$  using the residuals of the  
 167  $(100-X)$ -th percentile values for  $\mathcal{E}_{G,100-X}^\theta$  relative to the red regression line Figure 3A, which we  
 168 denote  $\mathcal{R}_{G,X}^\theta$ :

$$\mathcal{S}_{G,X}^\theta = \mathcal{R}_{G,X}^\theta / \widehat{\mathcal{D}_{G,X}^\theta} = (\mathcal{E}_{G,100-X}^\theta - \widehat{\mathcal{E}_{G,100-X}^\theta}) / \widehat{\mathcal{D}_{G,X}^\theta} = (\mathcal{D}_{G,X}^\theta - \widehat{\mathcal{D}_{G,X}^\theta}) / \widehat{\mathcal{D}_{G,X}^\theta}$$

169 where  $\widehat{\mathcal{D}_{G,X}^\theta}$  is the expected dispersion of dependency scores based on the robust linear  
 170 regression of  $\mathcal{E}_{G,100-X}^\theta$  given  $\mathcal{E}_{G,X}^\theta$ , or

$$\widehat{\mathcal{D}_{G,X}^\theta} = \widehat{\mathcal{E}_{G,100-X}^\theta} - \mathcal{E}_{G,X}^\theta.$$

172 The expected dispersion  $\widehat{\mathcal{D}_{G,X}^\theta}$  increases for more strongly negative efficacy scores  $\mathcal{E}_{G,X}^\theta$ ,  
 173 indicating greater variances in dependency effects for commonly or selectively essential genes  
 174 (e.g., greater variance for RAN vs. ZCWPW1 in Figure 3B). This could be a result of  
 175 experimental noise (e.g., fewer sequencing reads for negatively selected genes) or greater  
 176 biological variability in dependency effects for these genes. To better distinguish whether  
 177 genes have dispersion greater than would be expected simply based on their efficacy scores,  
 178 we normalize the residual values  $\mathcal{R}_{G,X}^\theta$  by dividing by the expected dispersion  $\widehat{\mathcal{D}_{G,X}^\theta}$  to obtain a  
 179 measure of selectivity that accounts for the variation of  $\widehat{\mathcal{D}_{G,X}^\theta}$  with the  $X$ -th percentile efficacy,  
 180  $\mathcal{E}_{G,X}^\theta$ . Both the efficacy and the selectivity of all the genes across different  $\theta$  is available to  
 181 download (Figure 3C – source data 3).

182

## 183 **Characteristics of essential genes**

184 By comparing the efficacy and the selectivity, we found that genes with strongly  
185 negative efficacy tend to be less selective, whereas selectively essential genes tend to have  
186 only moderate efficacy (Figure 3C). To characterize the genes that have either negative  
187 efficacy or positive selectivity, we performed gene set enrichment analysis of 6,551 pathways.  
188 This revealed the pathways overrepresented among essential, selective, and both selective  
189 and essential genes. For example, chromatin regulation is overrepresented among genes that  
190 are both selective and essential; nuclear metabolism and translation are overrepresented  
191 among essential (but not selective) genes; and regulation of kinase activity and tissue  
192 development are overrepresented among selective (but not essential) genes (Figure 3D,  
193 Figure 3D – source data 4).

## 194 195 **Relationship of the selectivity and the lineage specificity**

196 Cells in different lineages tend to depend on distinct essential genes compared to cells  
197 in the same lineage. We examined the extent to which lineage-specific dependence  
198 contributes to selectivity. For each gene, we assessed the relationship between the number of  
199 cell lines dependent on the gene and the gene's efficacy and selectivity, and confirmed that  
200 lower selectivity and more negative efficacy are associated with a greater number of  
201 dependent cell lines (Figure 4A).

202 We computed the number of distinct lineages dependent on each gene using the  
203 Adaptive Daisy Model (ADaM), a permutation-based statistical model reported previously  
204 (Behan et al., 2019). As with dependent cell lines, a greater number of dependent lineages  
205 was associated with more negative efficacy and lower selectivity (Figure 4B). Overall, we

206 found that 1,050 genes are commonly essential across all the lineages, 670 genes are  
207 essential in at least one lineage, and 2,581 essential genes were not lineage-dependent  
208 (Figure 4B-C, Figure 4 – source data 5). Unsurprisingly, the number of dependent cell lines is  
209 strongly associated with the number of dependent lineages (Figure 4C).

210 Using an independent CRISPR screening dataset, Behan et al. also proposed  
211 candidate drug targets based on selectively essential genes they identified for each lineage.  
212 We found that most genes with high selectivity scores were also identified as targets for one or  
213 more lineages in their analysis (Figure 4D). On the other hand, targets proposed for multiple  
214 lineages in Behan et al. (orange and red points in Figure 4D) tended to show moderate  
215 efficacy scores, but not necessarily high selectivity scores.

216 Though we saw a strong relationship between lower selectivity and a greater number of  
217 dependent cell lines and lineages (Figs. 4A-B), we note that lineage specificity is not the only  
218 cause of high selectivity. Conventionally, high selectivity is interpreted as being commonly  
219 essential within a few lineages but non-essential in others. Some genes do manifest this type  
220 of selectivity (e.g., broad dependence on CTNNB1 in colorectal cancer, MYB in leukemia, IRF4  
221 in multiple myeloma, KRAS in pancreatic cancer, and SOX10 in skin cancer, Figure 4E), but  
222 such common essentiality within a lineage is relatively unusual. In many more cases, a gene  
223 that shows high selectivity is selectively essential within each lineage (e.g., only partial  
224 dependence on CTNNB1 in liver, lung, and pancreatic cancers, or on IRF4 in skin cancer).  
225 CDK4 is particularly strong example of this pattern as it is not commonly essential in any  
226 lineages, but is selectively essential in many. While more negative efficacy  $\mathcal{E}_{G,X}^{\theta}$  is strongly  
227 correlated with larger dispersion  $\mathcal{D}_{G,X}^{\theta}$  (e.g., PSMA1 vs ISL1, Figure 4E; see also Figs. 3A-B),  
228 some commonly essential genes that show similar efficacy have substantially higher selectivity

229 than others (e.g., the selective PSMB5 vs. the non-selective PSMA1) (Figure 4E; see also Figs  
230 4A and 4C). These genes that are essential in many lineages but nevertheless have high  
231 selectivity could be promising drug targets given appropriate biomarkers to characterize  
232 sensitive cell lines.

233

234 **Clustering essential genes to find related targets**

235 The dependency profile of a gene carries information about the functions of the gene  
236 that make it essential in certain cellular contexts. When a set of genes comprise a functional  
237 unit (e.g., a pathway or a complex) that regulates cell viability, these genes would be expected  
238 to have similar dependency profiles. Gene-wise cluster analysis of the dependency data  
239 should therefore reveal functional units that connect essential genes into pathways or protein  
240 complexes. It may also be easier to relate the essentiality of pathways to cancer genotypes  
241 than to interpret the essentiality of individual genes. Pathway analysis can help identify  
242 druggable vulnerabilities at the pathway level that might be missed by single-gene analysis.

243 We clustered essential genes based on the similarity of the combined CRISPR-shRNA  
244 dependency scores across 423 cell lines. Our approach is based on a related pair of popular  
245 algorithms, t-distributed stochastic neighbor embedding (t-SNE) and density based spatial  
246 clustering and noise (DBSCAN) (Maaten and Hinton, 2008; Maaten, 2014; Ester et al., 1996).  
247 t-SNE is a technique that reduces the dimensionality of multi-dimensional data while  
248 preserving the pairwise distances between data points at high dimensions as much as  
249 possible. It has been widely used for visualizing high dimensional data, such as single-cell  
250 RNA-seq data (Mass et al., 2016). DBSCAN is a clustering algorithm that detects regions  
251 where the data points are gathered at high density and clusters them; it is often used to cluster

252 data points based on their coordinates in the t-SNE plot. The combination of t-SNE and  
253 DBSCAN (expressed as 't-SNE + DBSCAN', hereafter) is a powerful clustering algorithm for  
254 high-dimensional data, such as single-cell transcriptomes (Haber et al., 2017).

255 One limitation of this approach is that the t-SNE algorithm is stochastic, producing  
256 different results and clusters with different initial seeds. However, when we compared clusters  
257 yielded by t-SNE + DBSCAN from multiple runs, we found that strongly positively correlated  
258 points are always clustered together while weakly positively correlated points are less  
259 consistently so. To obtain robust cluster assignments from t-SNE + DBSCAN we therefore  
260 used a workflow we call ensemble clustering with hierarchy over DBSCAN on t-SNE with  
261 Spearman distance matrix (ECHODOTS).

262 Briefly, ECHODOTS consists of four steps (Figure 5A, Figure 5 – figure supplement 1):  
263 it 1) computes the pairwise Spearman distance matrix among essential genes, 2) feeds the  
264 distance matrix as input to run t-SNE with different initial seeds 200 times, 3) clusters data  
265 points based on their coordinates in the t-SNE plot with DBSCAN, and 4) identifies the sets of  
266 genes assigned to the same cluster consistently across the 200 sets of clusters using a  
267 technique called ensemble clustering (Hornik, 2005). ECHODOTS produces more reliable  
268 clusters than a single run of t-SNE + DBSCAN by seeking data points that are consistently  
269 clustered together.

270

### 271 **Cluster reveals known and new connections among essential genes**

272 We ran ECHODOTS against the combined dependency score  $S^\theta$  of the 4,301 essential  
273 genes ( $\theta=0.6$ ,  $X=1$ ), and assigned them into 879 small, 608 medium, and 338 large clusters  
274 (Figure 5 – source data 6). Genes in the same cluster tended to be close to each other on t-



275 SNE maps from individual runs, and some clusters were enriched for genes known to be  
276 members of specific biological pathways or complexes (Figure 5B, Figure 5B – source data 7).  
277 The median efficacy and selectivity of the clusters varied widely, suggesting that some  
278 represent more promising sets of drug targets than others (Figure 5C-D).

279 In examining the clusters, we found that they often included genes that were not all  
280 mutually correlated with one another. While strongly positively correlated genes tend to be  
281 located in the same neighborhood on the t-SNE map and subsequently clustered together in  
282 ECHODOTS (Figure 5B), for a gene to be added to a cluster it only needs to be correlated with  
283 at least one other gene in the cluster. The structure of the correlations among the genes within  
284 a cluster can therefore highlight subtle functional relationships. For example, we plotted the  
285 correlations among six highly essential genes in cluster S152, with an edge between genes  
286 indicating a Spearman correlation greater than 0.1 (Figure 5E). In this cluster, KEAP1 and  
287 KCTD10 are both strongly correlated with the E3 ubiquitin ligase CUL3 (with correlations of  
288 0.387 and 0.29, respectively), but have no correlation with each other (correlation -0.002). This  
289 is likely due to the fact that KEAP1 and KCTD10 interact with CUL3 in a mutually exclusive  
290 manner: both serve as adaptor proteins binding the same site on CUL3 but the resulting  
291 complexes degrade distinct target proteins (NFE2L2 and RHOB, respectively) (Cullinan et al.,  
292 2004; Kovačević et al., 2018).

293 Perhaps more intriguing are clusters that appear to show a connection between specific  
294 cellular processes and genes not otherwise known to be involved in that process. We offer two  
295 examples. One is Cluster L119 (Figure 5F), which is comprised of three small clusters.  
296 Clusters S369 and S745 contain the core MAP kinase (MAPK) pathway proteins, including  
297 KRAS, RAF1, BRAF, and MAPK1, while cluster S641 consists of CTNNB1 ( $\beta$ -catenin) and

TCF7L2, which form a bipartite transcription factor complex that is a key effector of the Wnt signaling pathway (Jin and Liu, 2008). These small clusters are within the same large cluster, suggesting that MAPK and Wnt signaling are functionally related in dependent cancer cell lines (Jeong et al., 2018). Cluster L119 also contains SHOC2, which was positively correlated to KRAS, RAF1, BRAF, and MAPK1. Multiple KRAS-mutant cancers were recently shown to be vulnerable to the loss of SHOC2 in the context of MEK inhibition, confirming the link between SHOC2 and MAP kinase pathway driven cancers (Sulahian et al., 2019).

A final example is the cluster L91 (Figure 5G), also consisting of multiple smaller clusters. Cluster S154 contains the selenoprotein GPX4. GPX4 encodes glutathione peroxidase 4, an antioxidant enzyme, that reduces cytotoxic lipid peroxides and protect cells from a non-apoptotic cell death, called ferroptosis. Intriguingly, we found the other genes in Cluster S154 were involved in selenoprotein synthesis (SEPHS2, SEPSECS, PSTK, EFFSEC) (Squires and Berry, 2008), suggesting that the primary role of these genes in dependent cell lines is to synthesize GPX4. S572 contains another selenoprotein TXNRD1 and its substrate, TXNDC17 (Espinosa and Arnér, 2018), both of which are also strongly correlated with four selenoprotein synthesis genes. Overall, L91 seems to represent a gene set related to the sensitivity to ferroptosis (Abdalkader et al., 2018; Ingold et al., 2018).

315

## 316 **Comparison between different dependency scores**

We have so far computed the efficacy, the selectivity, and the clusters of essential genes using  $S^\theta$  with the fixed mixing ratio,  $\theta=0.6$ , since this  $S^\theta$  retains the largest variance of the original CRISPR and shRNA scores (Figure 2A). Here we compare the performance of different  $\theta$ .

320

321 For efficacy, larger  $\theta$  (i.e., with a larger contribution of  $S^C$  to  $S^\theta$ ) gave more genes with  
322 strongly negative efficacy (Figure 5 – figure supplement 2A). Consistently, the number of  
323 dependent cell lines and the number of dependent lineages per gene increased with larger  $\theta$   
324 (Figure 5 – figure supplement 2B-F). For selectivity, larger  $\theta$  gave more genes with high  
325 selectivity (Figure 5 – figure supplement 2G). We showed earlier that essential (i.e., negative  
326 efficacy), selective (i.e., high selectivity), and both selective and essential genes overrepresent  
327 different pathways when  $\theta = 0.6$  (Figure 3D). Similar pathways were overrepresented by  
328 selective and both selective and essential genes when  $\theta = 0.8$  and 1. However, no pathways  
329 were associated with selective or both selective and essential genes when  $\theta < 0.5$ , suggesting  
330 that a high selectivity was given to genes more randomly (Figure 5 – figure supplement 2H).  
331 Since it is more likely that selective genes represent certain pathways such as the ones shown  
332 in Figure 2D,  $S^\theta$  with  $\theta > 0.5$  are more reasonable ones to choose.

333 Next, we compared the clusters of the essential genes using ECHODOTS. Since the  
334 number of discovered essential genes varies with different  $\theta$  (Figure 2D), we expect the  
335 number of clusters to be different. Therefore, we did not fix the number of clusters across  $\theta$ .  
336 Instead, we sought the upper bound of the neighborhood threshold  $\varepsilon$  (termed  $\varepsilon_0$ ) in DBSCAN  
337 for each  $\theta$  because as  $\varepsilon$  gets larger than a certain value,  $\varepsilon_0$ , most points on the t-SNE plot  
338 would start to merge to form a single large cluster (See Methods). We can detect incorrect  
339 merging by measuring the ratio between the 1<sup>st</sup> and 2<sup>nd</sup> largest cluster sizes (Figure 5 – figure  
340 supplement 3A-B). We found that  $\varepsilon_0$  is particularly small for  $\theta = 0.8$  and 1 compared to the rest  
341 of  $\theta$ , and more clusters were discovered for these  $\theta$  consequently (Figure 5 – figure  
342 supplement 3C-D). We compared cluster memberships of the 2,008 genes identified as  
343 essential among all  $\theta$ , and found that  $\theta > 0.5$  and  $\theta < 0.5$  gave substantially different clusters

(Figure 5 – figure supplement 3E). Through the comparison, we concluded that our initial choice of  $\theta = 0.6$  was a reasonable one since the combined dependency score is more informative with more weight on CRISPR than shRNA ( $\theta > 0.5$ , Figure 5 – figure supplement 3) while having some contribution from shRNA is more beneficial than CRISPR alone (Figure 2E).

### **shinyDepMap: an interactive web tool to explore the essentiality of genes**

Both the clusters of essential genes and the gene efficacy and selectivity scores provide valuable information for finding potential chemotherapeutic drug targets. To make this information accessible to the broader community of experimental drug discovery researchers, we developed a web-based tool to explore these analyses, called shinyDepMap. shinyDepMap is written in R (Chang et al., 2019) using the shiny package for building interactive visualization tools. It consists of two apps: “Gene essentiality” and “Gene cluster.” Each app is a dashboard-style website (Figure 6A). shinyDepMap can be used in three ways: 1) via the website <https://labsyspharm.shinyapps.io/depmap>, 2) by downloading the code and pre-processed data from the GitHub repository (<https://github.com/kenichi-shimada/shinyDepMap>) and running it on a local computer, and 3) running the app from a Docker container using the image at <https://hub.docker.com/r/labsyspharm/shinydepmap>. The analysis workflow in the application is explained below (Figure 6B).

**Gene essentiality (all protein-encoding genes)** This app allows a user to explore the essentiality of all the genes tested in the DepMap genetic perturbation experiments. Its output has two panels. A scatterplot in the middle displays efficacy and selectivity scores for all genes (3, bold numbers correspond to the panels in Figure 6B-C). By hovering over the plot points

367 with the cursor, one can find the genes corresponding to each point. When a gene name to  
368 search is provided in the input text box (1), corresponding genes will be highlighted in  
369 orange/red in the Efficacy-Selectivity plot (3). Genes matched with the query will be listed on  
370 the “Matched genes” tab in the right (4), in which by clicking a gene’s name, the description of a  
371 gene in GeneCards (<https://www.genecards.org/>) will be open on a new page. By further  
372 selecting a matched gene from the dropdown menu (2), one can see the combined  
373 dependency scores of the gene in 423 cell lines in the “Dependency scores” tab on the right  
374 (5). The definition of combined dependency scores, the efficacy, and the selectivity can be  
375 changed by tuning the mix ratio (equivalent to  $\theta$ ; 6) and the efficacy threshold (X-th percentile,  
376 7) from the input panel. We set them  $\theta=0.6$  and  $X=1$  by default.

377 **Gene cluster (essential genes)** This app allows a user to explore gene clusters among  
378 the essential genes. When a user first selects an essential gene from the top-left dropdown  
379 menu (8), genes clustered with the query gene will be shown on the output panels. There are  
380 three output panels in the app. The top-center is the Efficacy-Selectivity plot for the essential  
381 genes (9). The bottom-center shows the t-SNE plot, indicating the similarity of the dependency  
382 scores among essential genes (10). The list of “Clustered genes” will be shown in the right  
383 panel (11). “Connectivity” tab will show how the clustered genes are connected (i.e., strongly  
384 correlated) (12). The graphs can be downloaded in the GraphML format. “Correlation” tab  
385 shows the Spearman correlation coefficients between the selected gene and the other  
386 essential genes grouped by the clusters (13). While all the genes in the Small cluster are  
387 shown by default, one can change it by tuning the “cluster size” parameter (14) and the  
388 probability threshold (15) in the left input panel. In ECHODOTS, we computed a probability at  
389 which each gene belongs to the assigned cluster. By setting the probability threshold close to

one, one can show only genes that are assigned to the same clusters consistently across many runs of t-SNE + DBSCAN. This app allows the users to tune the mix ratio (**16**) and the efficacy threshold (**17**) like the Gene essentiality app. Essential genes were defined based on these two parameters. Consequently, modifying these parameters affects the clusters.

## Discussion

In this paper we described an interactive software tool, shinyDepMap, that allows users to rapidly determine the efficacy and selectivity of a gene of interest and thereby find highly selective genes that may offer promising therapeutic targets. shinyDepMap is based on both the CRISPR and shRNA genome-wide screening DepMap datasets, which we combined to generate a unified dependency score that is more informative than data from either dataset alone. Using this combined dependency score we computed “efficacy” and “selectivity” scores for each gene and highlighted how these scores can be used to characterize the therapeutic potential of targeting different genes across cell lines and lineages. Finally, we performed robust clustering of commonly and selectively essential genes to highlight functional relationships. shinyDepMap allows users to interactively explore both the essentiality and clustering results and is available as a deployed web-based application at <https://labsyspharm.shinyapps.io/depmap> and via source code or Docker image.

Our cluster analysis of the dependency scores highlighted genes comprising complexes and pathways, as previously reported (Pan et al., 2018). Our research complements published work in the area of complex and pathway annotation, in part because we were able to combine both DepMap datasets. We provide cluster information in a browsable form in shinyDepMap, including the ability to tune the size of clusters. One application of this tool is “target hopping,”

*i.e.*, moving from one drug target to another while keeping the selectivity profile similar (Schenone et al., 2013). One goal of target hopping is to identify druggable targets with similar dependency score profiles to genes of interest that are not conventionally druggable. A classic example is KRAS, which is essential to many cancers but until recently has been considered “undruggable.” KRAS appears in cluster L119 with the druggable kinases RAF1 and MAPK1, highlighting these proteins as relevant alternative targets. One goal of the shinyDepMap is to help researchers identify similar therapeutic opportunities among less well-studied genes.

Despite its potential value for therapeutic discovery, the DepMap dataset and our corresponding analysis in shinyDepMap has important practical limitations. The DepMap data characterizes the genetic requirements for cells grown in culture, which differs from the *in vivo* tumor environment in critical ways: the presence of nutrient-rich media, a two-dimensional rather than a three-dimensional substrate, and the lack of a functional immune system or physiological microenvironment. In addition, the effects of genetic perturbations (knockdown and knockout) do not necessarily correspond to those of chemical inhibition, a discrepancy that makes target identification from datasets like the DepMap less straightforward (Weiss et al., 2007). Finally, because the DepMap includes only cultured cancer lines and no wild-type cell lines or tissues, a highly selective gene in our analysis is not guaranteed to be less toxic to normal tissues when inhibited. Provided these limitations are kept in mind, the DepMap is a powerful resource and we hope the shinyDepMap tool makes it accessible to a broad community of researchers.

## **Methods**

### ***Data and Code Availability Statement***

436           Following data of the 2019 Q3 release were downloaded from the DepMap project  
437 website: CRISPR (avana) ("Achilles\_gene\_effect.csv"), combined RNAi  
438 ("D2\_combined\_gene\_dep\_scores.csv"), and the cell line metadata ("sample\_info.csv"). The  
439 CRISPR and shRNA efficacy data provided by the DepMap project were normalized with  
440 CERES and DEMETER2 algorithms by the Broad Institute, respectively. To compute the  
441 combined dependency score, we use the data of 15,847 genes in 423 cell lines, which were  
442 examined with both CRISPR and shRNA. The codes generated during this study are available  
443 at <https://github.com/kenichi-shimada/depmap-analysis> (data processing and analysis) and  
444 <https://github.com/kenichi-shimada/shinyDepMap> (standalone shinyDepMap). This study did  
445 not generate unique datasets.

446

#### 447 ***Imputing missing values in shRNA and CRISPR dependency scores***

448           We first took conditions that were tested in both CRISPR and shRNA were plotted.  
449 15,847 genes were perturbed using both methods in 423 cell lines, and we compared  
450 4,846,055 conditions that were non-missing values. Here, one condition is defined as a  
451 perturbation of one gene, either using CRISPR or shRNA in one cell line. We next imputed  
452 missing values in CRISPR and shRNA datasets using non-missing data from the other  
453 method. 1,345,642 conditions (20%) were tested using CRISPR, but not shRNA. 19,001  
454 (0.28%) were tested using shRNA, but not CRISPR. In these cases, the missing values were  
455 imputed from the other data using local polynomial regression (loess) function in R. 4,457  
456 (0.066%) conditions were not tested by either CRISPR or shRNA, which were left as missing.  
457 After missing values were imputed, the combined dependency score from two values was  
458 computed.



459

### 460 ***Computing the dependency score combining CRISPR and shRNA data***

461 In this method section, we use  $S_{G,L}^\theta$  instead of  $S^\theta$  to represent the dependency score of  
462 gene  $G$  in cell line  $L$ , to make the argument clearer. To summarize CRISPR and shRNA  
463 dependency scores, we developed a new dependency score of a gene  $G$  in a cell line  $L$ ,  $S_{G,L}^\theta$ ,  
464 using the following equation:

$$S_{G,L}^\theta = \theta S_{G,L}^C + (1 - \theta) S_{G,L}^R \quad (0 \leq \theta \leq 1)$$

465 where  $S_{G,L}^C$  and  $S_{G,L}^R$  are the dependency scores of the gene  $g$  in a cell line  $l$ , given by CRISPR  
466 and shRNA alone, respectively. The combined dependency score is a function of the mixing  
467 value  $\theta$ . We particularly chose and compared six values of  $\theta: \theta \in \{0, 0.2, 0.4, 0.6, 0.8, 1\}$ .  
468 Resulting combined dependency scores,  $S_{G,L}^\theta$ , were computed for 15,847 genes in 423 cell  
469 lines for each  $\theta$ .

470

### 471 ***Defining essential conditions***

472 We defined the loss of a gene  $G$  is essential in a cell line  $L$  when the score  $S_{G,L}^\theta$  is lower  
473 than the essentiality threshold  $T_\theta$ :  $S_{G,L}^\theta < T_\theta$ .  $T_\theta$  is defined as follows: we first fit the kernel  
474 density estimate function to the entire distribution of the dependency scores of across all  
475 genes and cell lines, or  $S_{:,L}^\theta$ . This distribution is well fitted with a normal distribution with a heavy  
476 left tail. We computed the mean  $\mu$  and standard deviation  $\sigma$  from the right-half of the data and  
477 defined  $T_\theta$  such that  $P(S_{G,L}^\theta < T_\theta) = 0.001$  where  $S_{G,L}^\theta \sim N(\mu, \sigma)$ .

478

### 479 ***Identifying inconsistent essential genes between CRISPR and shRNA***

480 To find out genes that are claimed essential only by CRISPR or shRNA, we first  
 481 computed  $T_\theta$  for CRISPR and shRNA, or  $T_C$  and  $T_R$ . Next, we sought the set of scores  
 482 targeting the same gene in all the cells ( $S_{G,\cdot}^C, S_{G,\cdot}^R$ ). The essentiality claimed only by CRISPR or  
 483 shRNA were expressed as  $S_{G,\cdot}^C < T_C \cap S_{G,\cdot}^R \geq T_R$  and  $S_{G,\cdot}^C \geq T_C \cap S_{G,\cdot}^R < T_R$ , respectively. They were  
 484 illustrated as areas A and B in Fib. 1B. Using one-tailed Fisher's exact test, we computed the  
 485 statistical significance of the enrichment of the data points ( $S_{G,\cdot}^C, S_{G,\cdot}^R$ ) in the areas A and B and  
 486 found 958 and 20 genes were claimed essential only by CRISPR and shRNA, respectively. We  
 487 also computed the statistical significance of the overlap between these gene sets with publicly  
 488 available gene annotations, Molecular Signature Database (MSigDB v7.0) (Subramanian et al.,  
 489 2005), using Fisher's exact test.

490

### 491 ***Computing the efficacy and selectivity for each gene***

492 The efficacy  $\mathcal{E}_{G,X}^\theta$  measures how essential a gene is in a sensitive cell line:

$$493 \quad \mathcal{E}_{G,X}^\theta = \text{the } X^{\text{th}} \text{ percentile of } S_{G,\cdot}^\theta \quad (X \in \{1, 2.5, 5, 10, 25\})$$

494 We defined a gene as essential for a given  $X$  when  $\mathcal{E}_{G,X}^\theta < T_\theta$  (Figure 3A). Since  $T_\theta$  is selected  
 495 such that the probability  $P(S_{G,L}^\theta < T_\theta) = 0.001$  when  $G$  is non-essential, we computed the  
 496 probability of  $N$  cell lines of

497 To define the selectivity, we first determined the dispersion of the distribution of  $S_{G,\cdot}^C$ ,  
 498  $\mathcal{D}_{G,X}^\theta$ :

$$\mathcal{D}_{G,X}^\theta = \mathcal{E}_{G,100-X}^\theta - \mathcal{E}_{G,X}^\theta$$

499  $\varepsilon_{G,X}^\theta$  and  $\varepsilon_{G,100-X}^\theta$  were in a strong linear relationship for most of the genes, which correspond  
 500 to commonly essential genes, while some genes have large positive residuals. We defined the  
 501 residual from the regression line,  $\mathcal{R}_{G,X}^\theta$ , as follows:

$$\mathcal{R}_{G,X}^\theta = \mathcal{D}_{G,X}^\theta - \widehat{\mathcal{D}_{G,X}^\theta}$$

$$\widehat{\mathcal{D}_{G,X}^\theta} = \widehat{\varepsilon_{G,100-X}^\theta} - \varepsilon_{G,X}^\theta = f(\varepsilon_{G,X}^\theta) - \varepsilon_{G,X}^\theta$$

502 and the selectivity of the gene  $G$ ,  $\mathcal{S}_{G,X}^\theta$ , was defined as:

$$\mathcal{S}_{G,X}^\theta = \mathcal{R}_{G,X}^\theta / \widehat{\mathcal{D}_{G,X}^\theta}.$$

504

### 505 **Overlap of essential genes among different $\theta$**

506 The number of essential genes, defined as  $\varepsilon_{G,1}^\theta < T_\theta$ , depends on the mixing ratio  $\theta$ . To  
 507 assess the overlap between the essential gene sets across different  $\theta$ , we computed an  
 508 overlap index for any pairs of  $\theta$  that is denoted as  $O_X(\theta_1, \theta_2)$ :

$$O_X(\theta_1, \theta_2) = N_X^{\theta_1 \cap \theta_2} / \min(N_X^{\theta_1}, N_X^{\theta_2})$$

509 where  $\theta_1$  and  $\theta_2$  are specific values of  $\theta$ ,  $N_X^{\theta_1}$  is the number of essential genes when  $\theta = \theta_1$ ,  
 510 and  $N_X^{\theta_1 \cap \theta_2}$  is the number of shared genes between the two essential gene sets. By definition,  
 511  $O_X(\theta_1, \theta_2)$  can take any values between 0 and 1: the index is zero when the two essential gene  
 512 sets do not share any genes; the index is one when the two essential genes are identical  
 513 (Figure 2D).

514

### 515 **Identifying pathways overrepresented by essential and/or selective genes**

516 For each of essential (i.e., negative efficacy), selective (i.e., high selectivity), or both  
 517 selective and essential genes, we sought pathways that were overrepresented by each gene.

518 We sorted all the genes by the efficacy and the selectivity in descending order and ran gene  
519 set enrichment analysis (GSEA) with the sorted genes against the pathways from MSigDB.  
520 GSEA was performed utilizing fgsea package with  $10^7$  permutations (Sergushichev, 2016).

521

### 522 ***Identifying lineage-specific and universally essential genes***

523 To compute the lineage specificity, we utilized the Adaptive Daisy Model (ADaM). ADaM  
524 calculates the minimum number of dependent cell lines that are required for a gene to be  
525 considered as commonly essential among the cell lines in question (Behan et al., 2019). It is  
526 implemented in the ADAM2 R package (<https://github.com/DepMap-Analytics/ADAM2>). The  
527 dependency score matrix,  $S_{:,j}^{\theta}$ , contains the information of 423 cell lines representing 28  
528 lineages. We focused on a subset of 387 cell lines in 17 lineages that includes ten or more cell  
529 lines. We computed the binary essentiality matrix for each G and L, where 1 if  $S_{G,L}^{\theta} < T_{\theta}$  and 0  
530 otherwise. We then provided a subset of the matrix attributed to each lineage as input and  
531 calculated the minimum number of dependent cell lines for the lineage. Each lineage is  
532 considered dependent on G when the number of dependent cell lines in the lineage is equal to  
533 or greater than the minimum number of the dependent cell lines. To compute universally  
534 essential genes utilizing ADaM, we calculated the binary essential matrix for the 17 lineages  
535 instead of cell lines and the minimum number of dependent lineages providing the matrix as  
536 input.

537 Behan et al. also provides a list of genes that are good targets for chemotherapies for  
538 each lineage. We counted the number of lineages they suggested was a good target for each  
539 gene and mapped them onto the Efficacy/Selectivity plot. (Figure 4D).

540

## 541 ***Robust cluster analysis utilizing t-SNE and DBSCAN: ECHODOTS***

542 We implemented a new cluster algorithm, ensemble clustering with hierarchy over  
543 DBSCAN on t-SNE with Spearman distance matrix (ECHODOTS), extending the combination  
544 of t-SNE and DBSCAN. It is graphically summarized in Figure 5A and its pseudocode is  
545 provided in Figure 5 – figure supplement 1.

546 First, we computed the Spearman distance matrix or 1-Spearman correlation coefficient  
547 across all pairs of essential genes (line (1) in Figure 5 – figure supplement 1). This matrix was  
548 provided as input, and the coordinates of each data point in a 2D plane was computed with t-  
549 SNE (line (2)). Next, we clustered data points based on their coordinates with DBSCAN such  
550 that any two points whose distance is smaller than the neighborhood threshold  $\varepsilon$  are assigned  
551 into the same cluster (line (4)). We note that the range of the coordinates,  $L$ , varies among  
552 different runs of t-SNE. It is more reasonable to make the denominator  $d$  constant rather than  
553 to make  $\varepsilon$  constant, therefore we determined  $\varepsilon$  relative to  $L$ :  $\varepsilon = L/d$  (line (3)). The resulting  
554 cluster set  $C$  was computed by DBSCAN using  $\varepsilon$  derived from the preset parameter  $d$  (line  
555 (4)), which we set an integer ranging from 30 to 200. We ran t-SNE and DBSCAN to compute  
556  $C$  200 times using different initial seeds.

557 Next, we computed consistent clusters  $CC$  across 200 cluster sets  $C$  for each  $d$  (or  
558 equivalently  $\varepsilon$ ) using ensemble clustering available in clue R package (line (5)). When  $d$  is too  
559 small, ( $\varepsilon$  is too large), most of the points are erroneously connected. The mean size and the  
560 total number of consistent clusters depend on the neighborhood threshold  $\varepsilon$ . The smallest  $\varepsilon$   
561 yields the largest number of small clusters, in which data points within each cluster are most  
562 tightly connected. The largest  $\varepsilon$  yields the smallest number of large clusters, in which data  
563 points within each cluster are most loosely connected. When  $\varepsilon$  is too large (correspondingly,  $d$

was too small), most genes form one massive cluster in an extreme case. To avoid this, we set a lower bound for  $d$  or equivalently upper bound for  $\varepsilon$  such that  $\varepsilon = L/d \leq L/d_0 = \varepsilon_0$  (line (6)). We determined  $d_0$  by looking at the ratio between the sizes of the first and second largest clusters (Figure 5 – figure supplement 3A). When  $\varepsilon$  gets smaller ( $d$  gets larger), clusters get smaller and tighter. With smaller  $\varepsilon$ , some genes are not clustered with any other genes. We called these genes forming clusters only by themselves ‘noise’ and distinguished them from other clusters containing more than one gene. Eventually, all the genes become isolated or noise, but we stopped our analysis far ahead ( $d = 200$ ). When  $\theta = 0.6$ , we chose three different values of  $d$  ( $d = 65, 100, 141$ ), which discovered 338, 608, and 879 clusters (line (7)). Most of the smaller clusters derived with larger  $d$  were contained in the larger clusters derived with smaller  $d$ . Therefore, we constructed a hierarchical relationship between the three sets of clusters. There were a few cases where a gene belongs to distinct clusters with different  $d$ , but the smaller cluster is not a part of the larger cluster and thus the hierarchy is not constructed.

Note that DBSCAN is a hard clustering algorithm that assigns each gene into only one cluster. Once clustered, both strongly and weakly correlated genes become indistinguishable in hard clusters. On the other hand, ECHODOTS performs a soft clustering that can assign each gene to more than one cluster. It conducts the ‘majority vote’ among 200 runs of the clustering and computes the probability of a point being assigned to each cluster, which provides us with the strength of evidence for the cluster assignment of each gene. Thus, ECHODOTS produces more reliable clusters than a single run of t-SNE + DBSCAN by seeking data points that are consistently clustered together.

## ***Implementation of shinyDepMap website***

shinyDepMap was built using shiny package. Following packages are also used to implement the tool: ggplot2 (v3.2.1), RColorBrewer (v1.1-2), shinyWidgets (v0.4.8), plotly (v4.9.0), DT (v0.8), visNetwork (v2.0.8), tibble (v2.1.3), dplyr (v0.8.3), tidyr (v0.8.3). shinyDepMap can be run locally without an internet connection. One can download the code and data at <https://github.com/kenichi-shimada/shinyDepMap>, and run locally following the link's instruction.

## **Acknowledgments**

The authors thank Laura Maliszewski, Peter Sorger, and Rebecca Ward of Harvard Medical School for helpful discussions. This work is financially supported by the Japan Society for the Promotion of Science Overseas Research Fellowship (to K.S.), NIH R35GM131753 (to T.J.M.), and P50GM107618.

## **Competing Interests**

J.A.B. has received consulting fees from Two Six Labs, LLC. The other authors do not have competing interests to declare.

## 603 **Figure Legends**

### 604 **Figure 1. Systematic biases in CRISPR and shRNA dependency scores**

605 **A** Comparison of normalized CRISPR and shRNA dependency scores of 15,847 protein-  
606 encoding genes in 423 cell lines. The density and contour plot corresponds to the distribution  
607 of the scores from all gene perturbations. Vertical and horizontal solid lines indicate the  
608 essentiality thresholds for CRISPR and shRNA dependency scores, respectively. Areas A, B,  
609 C correspond to the regions where only CRISPR, only shRNA, or both CRISPR and shRNA  
610 claimed essential. **B** Comparison of CRISPR and shRNA targeting four genes. In each panel,  
611 data points correspond to each gene's perturbation in 423 cell lines. Each point corresponds to  
612 one cell line. **C** 958 and 20 genes were claimed essential only by CRISPR or shRNA but not  
613 by the other method, respectively (Fisher's exact test, p-value < 1e-3). **D** Assessment of the  
614 pathways overrepresented by the essential genes claimed only by CRISPR or shRNA,  
615 highlighted in **C** (Fisher's exact test).

616

### 617 **Figure 2. Identification of essential genes based on combined dependency score**

618 **A** Dependency scores defined with different mixing ratios are computed by projecting each  
619 point onto the corresponding lines.  $\theta$  denotes the fraction of CRISPR dependency scores. PC1  
620 is the direction of the primary principal component line. **B** The distributions of combined  
621 dependency scores for four genes showed in Figure 1B. **C** Top panel: the distribution of the  
622 combined dependency score  $S^{0.6}$ . The essentiality threshold  $T^{0.6}$  is determined based on this  
623 distribution. Bottom panel: the distribution of efficacy scores  $\mathcal{E}_{G,X}^{0.6}$  with various X (-th percentile).  
624 Genes that satisfy  $\mathcal{E}_{G,X}^{0.6} < T^{0.6}$  are defined as commonly or selectively essential, and the  
625 number of essential genes depends on X. **D** On the diagonal line, the numbers of commonly or



selectively essential genes identified with various  $\theta$  are shown. In the off-diagonal area, the numbers of essential genes identified with two distinct  $\theta$  are shown. Color code indicates an overlap index, a measure of overlap between two essential gene sets. The overlap index ranges from 0 (no shared genes) to 1 (the smaller set is included in the larger set). **E** The extent to which the genes claimed essential by only CRISPR or shRNA are covered by the essential genes discovered by each mixing ratio.

632

### 633 **Figure 3. The efficacy and the selectivity**

634 **A** The 1<sup>st</sup> and 99<sup>th</sup> percentiles of the combined dependency score of each gene where  $\theta = 0.6$ . Each point corresponds to one gene. X- and Y-axis are equivalent to the efficacy with X=1 and X=99, respectively. Solid red line and Dashed black line are robust linear regression and identity lines, respectively. **B** Distribution of the combined dependency scores of four selective and four non-selective genes. The 1<sup>st</sup> and 99<sup>th</sup> percentile values within each distribution are also highlighted. **C** The efficacy and the selectivity of all genes are plotted. **D** Summary of the pathways overrepresented by genes with strongly negative efficacy and high selectivity, strongly negative efficacy, and high selectivity, respectively.

642

### 643 **Figure 4. The lineage dependency**

644 In the following panels **A**, **B**, and **D**, the efficacy and the selectivity scatterplots (Figure 3C) are color-coded differently, highlighting the following properties of each gene. **A** The numbers of dependent cell lines. **B** The number of dependent lineages, computed with ADaM. **C** The relationship between the dependent cell lines and the dependent lineages. **D** The number of lineages in which Behan et al. suggested suitable for chemotherapy targets. **E** Nine genes'

649 dependency scores grouped by lineages, together with the number of dependent lineages and  
650 cell lines, the efficacy, and the selectivity. All the panels in this figure were computed using  $\theta =$   
651 0.6 and  $X=1$ .

652

### 653 **Figure 5. Essential gene clustering**

654 **A** The framework of ECHODOTS algorithm. **B** Nine gene clusters and their associated  
655 pathways. **C** Median efficacy and selectivity of Large Clusters. **D** Genes consisting Large  
656 Clusters with high selectivity highlighted in **C**. **E-G** The intra-cluster connectivity of three gene  
657 clusters as exemplars. The colors of nodes indicate their membership of small clusters, and  
658 the edges indicate that the two connected genes have Spearman correlation coefficient greater  
659 than 0.1. Numbers in **E** indicate Spearman correlation coefficients.

660

### 661 **Figure 6. shinyDepMap: a web-tool to explore DepMap dataset**

662 **A** Top pag. **B** Gene essentiality app. **1** Textbox to type in a (partial) gene symbol to query. **2**  
663 Dropdown menu to select a gene symbol that matches the query. **3** Efficacy-selectivity scatter  
664 plot. **4** List of matched genes. **5** Combined dependency score profile of the gene selected in **2**.  
665 **6** Mix ratio **7** Efficacy threshold

666 **C** Gene clusters app. **8** Dropdown menu to select an essential gene to explore. **9** Efficacy-  
667 selectivity plot for essential genes. **10** t-SNE plot. **11** List of clustered genes. **12** Connectivity  
668 plot. **13** Spearman correlation between the selected gene and the other essential genes. **14**  
669 Cluster size input. **15** Probability threshold input. **16** Mix ratio. **17** Efficacy threshold.

670

### 671 **Figure 5 – figure supplement 1. ECHODOTS algorithm**

672 ECHODOTS algorithm is written using pseudocodes. The line numbers correspond to the line  
673 numbers in the main text.

674

675 **Figure 5 – figure supplement 2. Efficacy, selectivity and the dependent lineages with**  
676 **various  $\theta$**

677 **A** Empirical cumulative density functions (CDF) of the efficacy  $\mathcal{E}_{G,X}^{\theta}$  across all the genes with  
678 various  $\theta$ . **B** Empirical CDFs of the number of dependent cell lines across all the genes with  
679 various  $\theta$ . **C** The distribution of the number of dependent lineages among essential genes with  
680 various  $\theta$ . The genes in the left- and right-hand side of the vertical lines are considered  
681 selectively and commonly essential according to ADaM. **D-E** The efficacy-selectivity plot of all  
682 the genes with various  $\theta$  ( $X=1$ ). The genes are color-coded based on the number of dependent  
683 cell lines (**D**) and lineages (**E**). **F** Relationship between the number of dependent lineages and  
684 the number of dependent cell lines with various  $\theta$  ( $X=1$ ). **G** Empirical CDFs of the selectivity  
685  $\mathcal{S}_{G,X}^{\theta}$  across all the genes with various  $\theta$ . **H** The number of overrepresented pathways  
686 associated with genes with strongly negative efficacy and high selectivity, strongly negative  
687 efficacy, and high selectivity.

688

689 **Figure 5 – figure supplement 3. Dependent cell lines and lineages using six dependency**  
690 **scores**

691 In the panels **A-D**, four parameters were plotted on Y-axis against  $d = L/\varepsilon$  on X-axis, where  $\varepsilon$   
692 is a neighborhood threshold in DBSCAN, for various for various  $\theta$ . **A** The ratio between the  
693 sizes of the first and second largest clusters ( $N_1/N_2$ ). **B** The number of genes assigned into the  
694 first and second largest clusters ( $N_1$  and  $N_2$ ) and the number of noise genes ( $N_n$ ), *i.e.*, the

695 genes which are not clustered with other genes. **C** The number of clusters. **D** The mean cluster  
696 size. **E** The similarity of the clusters with various  $\theta$ . Cluster membership of the 2,008 genes  
697 which were found essential with all  $\theta$  was compared using `cl_dissimilarity` in `clue` R package.

698

699

700

701 **Source Data**

702 Source Data were made publicly available from the following link:

703 [https://figshare.com/projects/shinyDepMap\\_Source\\_Data/97382](https://figshare.com/projects/shinyDepMap_Source_Data/97382)

704

705 **Figure 1 – source data 1. Information of the 423 cell lines in which both CRISPR and**  
706 **shRNA screening were tested.**

707 **Figure 2 – source data 2. The combined dependency scores for 15,847 protein-coding**  
708 **genes in 423 cell lines for six  $\theta$ .**

709  $\theta = 0$  and 1 corresponds to shRNA and CRISPR scores compared in Figure 1.

710 **Figure 3C – source data 3. Efficacy and Selectivity for 15,847 genes for the six  $\theta$  and five**

711 **X:  $X = \{1, 2.5, 5, 10, 25\}$**

712 **Figure 3D – source data 4. GO/KEGG Pathways overrepresented by genes with strongly**  
713 **negative efficacy, high selectivity, or strongly negative efficacy and high selectivity for**  
714 **six  $\theta$  and  $X=1$**

715 **Figure 4 – source data 5. Lineage-dependent essentiality of 17 lineages and common**  
716 **essentiality computed using ADaM for six  $\theta$**

717 **Figure 5 – Source Data 6. Cluster membership of essential genes and probability of their**  
718 **assignment to clusters for six  $\theta$**

719 **Figure 5B – Source Data 7. Pathways overrepresented in Large Clusters for six  $\theta$ .**

720 Clusters that contains 15 genes or more are only considered in this analysis.

721

722

723

724 **References**

- 725 Abdalkader M, Lampinen R, Kanninen KM, Malm TM, Liddell JR. 2018. Targeting Nrf2 to Suppress  
726 Ferroptosis and Mitochondrial Dysfunction in Neurodegeneration. *Front Neurosci* **12**.  
727 doi:10.3389/fnins.2018.00466
- 728 Bailey MH, Tokheim C, Porta-Pardo E, Sengupta S, Bertrand D, Weerasinghe A, Colaprico A, Wendl  
729 MC, Kim J, Reardon B, Ng PK-S, Jeong KJ, Cao S, Wang Zixing, Gao J, Gao Q, Wang F, Liu EM,  
730 Mularoni L, Rubio-Perez C, Nagarajan N, Cortés-Ciriano I, Zhou DC, Liang W-W, Hess JM,  
731 Yellapantula VD, Tamborero D, Gonzalez-Perez A, Suphavitai C, Ko JY, Khurana E, Park PJ,  
732 Van Allen EM, Liang H, Caesar-Johnson SJ, Demchok JA, Felau I, Kasapi M, Ferguson ML,  
733 Hutter CM, Sofia HJ, Tarnuzzer R, Wang Zhining, Yang L, Zenklusen JC, Zhang J (Julia),  
734 Chudamani S, Liu J, Lolla L, Naresh R, Pihl T, Sun Q, Wan Y, Wu Y, Cho J, DeFreitas T, Frazer  
735 S, Gehlenborg N, Getz G, Heiman DI, Kim J, Lawrence MS, Lin P, Meier S, Noble MS, Saksena  
736 G, Voet D, Zhang Hailei, Bernard B, Chambwe N, Dhankani V, Knijnenburg T, Kramer R,  
737 Leinonen K, Liu Y, Miller M, Reynolds S, Shmulevich I, Thorsson V, Zhang W, Akbani R,  
738 Broom BM, Hegde AM, Ju Z, Kanchi RS, Korkut A, Li J, Liang H, Ling S, Liu W, Lu Y, Mills GB,  
739 Ng K-S, Rao A, Ryan M, Wang Jing, Weinstein JN, Zhang J, Abeshouse A, Armenia J,  
740 Chakravarty D, Chatila WK, de Bruijn I, Gao J, Gross BE, Heins ZJ, Kundra R, La K, Ladanyi M,  
741 Luna A, Nissan MG, Ochoa A, Phillips SM, Reznik E, Sanchez-Vega F, Sander C, Schultz N,  
742 Sheridan R, Sumer SO, Sun Y, Taylor BS, Wang Jioajiao, Zhang Hongxin, Anur P, Peto M,  
743 Spellman P, Benz C, Stuart JM, Wong CK, Yau C, Hayes DN, Parker JS, Wilkerson MD, Ally A,  
744 Balasundaram M, Bowlby R, Brooks D, Carlsen R, Chuah E, Dhalla N, Holt R, Jones SJM,  
745 Kasaian K, Lee D, Ma Y, Marra MA, Mayo M, Moore RA, Mungall AJ, Mungall K, Robertson  
746 AG, Sadeghi S, Schein JE, Sipahimalani P, Tam A, Thiessen N, Tse K, Wong T, Berger AC,  
747 Beroukhim R, Cherniack AD, Cibulskis C, Gabriel SB, Gao GF, Ha G, Meyerson M,  
748 Schumacher SE, Shih J, Kucherlapati MH, Kucherlapati RS, Baylin S, Cope L, Danilova L,  
749 Bootwalla MS, Lai PH, Maglinte DT, Van Den Berg DJ, Weisenberger DJ, Auman JT, Balu S,  
750 Bodenheimer T, Fan C, Hoadley KA, Hoyle AP, Jefferys SR, Jones CD, Meng S, Mieczkowski  
751 PA, Mose LE, Perou AH, Perou CM, Roach J, Shi Y, Simons JV, Skelly T, Soloway MG, Tan D,  
752 Veluvolu U, Fan H, Hinoue T, Laird PW, Shen H, Zhou W, Bellair M, Chang K, Covington K,  
753 Creighton CJ, Dinh H, Doddapaneni H, Donehower LA, Drummond J, Gibbs RA, Glenn R, Hale  
754 W, Han Y, Hu J, Korchina V, Lee S, Lewis L, Li W, Liu X, Morgan M, Morton D, Muzny D,  
755 Santibanez J, Sheth M, Shinbrot E, Wang L, Wang M, Wheeler DA, Xi L, Zhao F, Hess J,  
756 Appelbaum EL, Bailey M, Cordes MG, Ding L, Fronick CC, Fulton LA, Fulton RS, Kandoth C,  
757 Mardis ER, McLellan MD, Miller CA, Schmidt HK, Wilson RK, Crain D, Curley E, Gardner J,  
758 Lau K, Mallery D, Morris S, Paulauskis J, Penny R, Shelton C, Shelton T, Sherman M,  
759 Thompson E, Yena P, Bowen J, Gastier-Foster JM, Gerken M, Leraas KM, Lichtenberg TM,  
760 Ramirez NC, Wise L, Zmuda E, Corcoran N, Costello T, Hovens C, Carvalho AL, de Carvalho  
761 AC, Fregnani JH, Longatto-Filho A, Reis RM, Scapulatempo-Neto C, Silveira HCS, Vidal DO,  
762 Burnette A, Eschbacher J, Hermes B, Noss A, Singh R, Anderson ML, Castro PD, Ittmann M,  
763 Huntsman D, Kohl B, Le X, Thorp R, Andry C, Duffy ER, Lyadov V, Paklina O, Setdikova G,  
764 Shabunin A, Tavobilov M, McPherson C, Warnick R, Berkowitz R, Cramer D, Feltmate C,  
765 Horowitz N, Kibel A, Muto M, Raut CP, Malykh A, Barnholtz-Sloan JS, Barrett W, Devine K,  
766 Fulop J, Ostrom QT, Shimmel K, Wolinsky Y, Sloan AE, De Rose A, Giulianti F, Goodman M,  
767 Karlan BY, Hagedorn CH, Eckman J, Harr J, Myers J, Tucker K, Zach LA, Deyarmin B, Hu H,  
768 Kvecher L, Larson C, Mural RJ, Somiari S, Vicha A, Zelinka T, Bennett J, Iacocca M, Rabeno B,

769 Swanson P, Latour M, Lacombe L, Têtu B, Bergeron A, McGraw M, Staugaitis SM, Chabot J,  
770 Hibshoosh H, Sepulveda A, Su T, Wang T, Potapova O, Voronina O, Desjardins L, Mariani O,  
771 Roman-Roman S, Sastre X, Stern M-H, Cheng F, Signoretti S, Berchuck A, Bigner D, Lipp E,  
772 Marks J, McCall S, McLendon R, Secord A, Sharp A, Behera M, Brat DJ, Chen A, Delman K,  
773 Force S, Khuri F, Magliocca K, Maithel S, Olson JJ, Owonikoko T, Pickens A, Ramalingam S,  
774 Shin DM, Sica G, Van Meir EG, Zhang Hongzheng, Eijckenboom W, Gillis A, Korpershoek E,  
775 Looijenga L, Oosterhuis W, Stoop H, van Kessel KE, Zwarthoff EC, Calatuzzolo C, Cuppini L,  
776 Cuzzubbo S, DiMeco F, Finocchiaro G, Mattei L, Perin A, Pollo B, Chen C, Houck J,  
777 Lohavanichbutr P, Hartmann A, Stoehr C, Stoehr R, Taubert H, Wach S, Wullich B, Kycler W,  
778 Murawa D, Wiznerowicz M, Chung K, Edenfield WJ, Martin J, Baudin E, Bubley G, Bueno R,  
779 De Rienzo A, Richards WG, Kalkanis S, Mikkelsen T, Noushmehr H, Scarpaccia L, Girard N,  
780 Aymerich M, Campo E, Giné E, Guillermo AL, Van Bang N, Hanh PT, Phu BD, Tang Y, Colman  
781 H, Evason K, Dottino PR, Martignetti JA, Gabra H, Juhl H, Akeredolu T, Stepa S, Hoon D, Ahn  
782 K, Kang KJ, Beuschlein F, Breggia A, Birrer M, Bell D, Borad M, Bryce AH, Castle E, Chandan  
783 V, Cheville J, Copland JA, Farnell M, Flotte T, Giana N, Ho T, Kendrick M, Kocher J-P, Kopp K,  
784 Moser C, Nagorney D, O'Brien D, O'Neill BP, Patel T, Petersen G, Que F, Rivera M, Roberts L,  
785 Smallridge R, Smyrk T, Stanton M, Thompson RH, Torbenson M, Yang JD, Zhang L, Brimo F,  
786 Ajani JA, Gonzalez AMA, Behrens C, Bondaruk J, Broaddus R, Czerniak B, Esmaeli B,  
787 Fujimoto J, Gershenwald J, Guo C, Lazar AJ, Logothetis C, Meric-Bernstam F, Moran C,  
788 Ramondetta L, Rice D, Sood A, Tamboli P, Thompson T, Troncoso P, Tsao A, Wistuba I,  
789 Carter C, Haydu L, Hersey P, Jakrot V, Kakavand H, Kefford R, Lee K, Long G, Mann G, Quinn  
790 M, Saw R, Scolyer R, Shannon K, Spillane A, Stretch J, Synnott M, Thompson J, Wilmott J, Al-  
791 Ahmadie H, Chan TA, Ghossein R, Gopalan A, Levine DA, Reuter V, Singer S, Singh B, Tien  
792 NV, Broudy T, Mirsaidi C, Nair P, Drwiega P, Miller J, Smith J, Zaren H, Park J-W, Hung NP,  
793 Kebebew E, Linehan WM, Metwalli AR, Pacak K, Pinto PA, Schiffman M, Schmidt LS, Vocke  
794 CD, Wentzensen N, Worrell R, Yang H, Moncrieff M, Goparaju C, Melamed J, Pass H,  
795 Botnariuc N, Caraman I, Cernat M, Chemencedji I, Clipca A, Doruc S, Gorincioi G, Mura S,  
796 Pirtac M, Stancul I, Tcaciuc D, Albert M, Alexopoulou I, Arnaout A, Bartlett J, Engel J, Gilbert  
797 S, Parfitt J, Sekhon H, Thomas G, Rassl DM, Rintoul RC, Bifulco C, Tamakawa R, Urba W,  
798 Hayward N, Timmers H, Antenucci A, Facciolo F, Grazi G, Marino M, Merola R, de Krijger R,  
799 Gimenez-Roqueplo A-P, Piché A, Chevalier S, McKercher G, Birsoy K, Barnett G, Brewer C,  
800 Farver C, Naska T, Pennell NA, Raymond D, Schilero C, Smolenski K, Williams F, Morrison C,  
801 Borgia JA, Liptay MJ, Pool M, Seder CW, Junker K, Omberg L, Dinkin M, Manikhas G, Alvaro  
802 D, Bragazzi MC, Cardinale V, Carpino G, Gaudio E, Chesla D, Cottingham S, Dubina M,  
803 Moiseenko F, Dhanasekaran R, Becker K-F, Janssen K-P, Slotta-Huspenina J, Abdel-Rahman  
804 MH, Aziz D, Bell S, Cebulla CM, Davis A, Duell R, Elder JB, Hilty J, Kumar B, Lang J, Lehman  
805 NL, Mandt R, Nguyen P, Pilarski R, Rai K, Schoenfield L, Senecal K, Wakely P, Hansen P,  
806 Lechan R, Powers J, Tischler A, Grizzle WE, Sexton KC, Kastl A, Henderson J, Porten S,  
807 Waldmann J, Fassnacht M, Asa SL, Schadendorf D, Couce M, Graefen M, Huland H, Sauter G,  
808 Schlomm T, Simon R, Tennstedt P, Olabode O, Nelson M, Bathe O, Carroll PR, Chan JM,  
809 Disaia P, Glenn P, Kelley RK, Landen CN, Phillips J, Prados M, Simko J, Smith-McCune K,  
810 VandenBerg S, Roggin K, Fehrenbach A, Kendler A, Sifri S, Steele R, Jimeno A, Carey F,  
811 Forgie I, Mannelli M, Carney M, Hernandez B, Campos B, Herold-Mende C, Jungk C,  
812 Unterberg A, von Deimling A, Bossler A, Galbraith J, Jacobus L, Knudson M, Knutson T, Ma D,  
813 Milhem M, Sigmund R, Godwin AK, Madan R, Rosenthal HG, Adebamowo C, Adebamowo SN,  
814 Boussioutas A, Beer D, Giordano T, Mes-Masson A-M, Saad F, Bocklage T, Landrum L,

815 Mannel R, Moore K, Moxley K, Postier R, Walker J, Zuna R, Feldman M, Valdivieso F, Dhir R,  
 816 Luketich J, Pinero EMM, Quintero-Aguilo M, Carlotti CG, Dos Santos JS, Kemp R,  
 817 Sankarankuty A, Tirapelli D, Catto J, Agnew K, Swisher E, Creaney J, Robinson B, Shelley CS,  
 818 Godwin EM, Kendall S, Shipman C, Bradford C, Carey T, Haddad A, Moyer J, Peterson L,  
 819 Prince M, Rozek L, Wolf G, Bowman R, Fong KM, Yang I, Korst R, Rathmell WK, Fantacone-  
 820 Campbell JL, Hooke JA, Kovatich AJ, Shriver CD, DiPersio J, Drake B, Govindan R, Heath S,  
 821 Ley T, Van Tine B, Westervelt P, Rubin MA, Lee JI, Aredes ND, Mariamidze A, Lawrence MS,  
 822 Godzik A, Lopez-Bigas N, Stuart J, Wheeler D, Getz G, Chen K, Lazar AJ, Mills GB, Karchin R,  
 823 Ding L. 2018. Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell*  
 824 **173**:371-385.e18. doi:10.1016/j.cell.2018.02.060  
 825 Behan FM, Iorio F, Picco G, Gonçalves E, Beaver CM, Migliardi G, Santos R, Rao Y, Sassi F, Pinnelli M,  
 826 Ansari R, Harper S, Jackson DA, McRae R, Pooley R, Wilkinson P, Meer D van der, Dow D,  
 827 Buser-Doepner C, Bertotti A, Trusolino L, Stronach EA, Saez-Rodriguez J, Yusa K, Garnett  
 828 MJ. 2019. Prioritization of cancer therapeutic targets using CRISPR-Cas9 screens. *Nature*  
 829 **568**:511-516. doi:10.1038/s41586-019-1103-9  
 830 Chang W, Cheng J, Allaire JJ, Xie Y, McPherson J. 2019. shiny: Web Application Framework for R.  
 831 Cullinan SB, Gordan JD, Jin J, Harper JW, Diehl JA. 2004. The Keap1-BTB Protein Is an Adaptor That  
 832 Bridges Nrf2 to a Cul3-Based E3 Ligase: Oxidative Stress Sensing by a Cul3-Keap1 Ligase.  
 833 *Molecular and Cellular Biology* **24**:8477-8486. doi:10.1128/MCB.24.19.8477-8486.2004  
 834 Espinosa B, Arnér ESJ. 2018. Thioredoxin-related protein of 14 kDa as a modulator of redox  
 835 signalling pathways. *British Journal of Pharmacology* 544-553.  
 836 doi:10.1111/bph.14479@10.1111/(ISSN)1476-  
 837 5381.nitric\_oxide\_the\_1998\_prize\_virtual\_issue.html  
 838 Ester M, Kriegel H-P, Sander J, Xu X. 1996. A Density-based Algorithm for Discovering Clusters a  
 839 Density-based Algorithm for Discovering Clusters in Large Spatial Databases with  
 840 NoiseProceedings of the Second International Conference on Knowledge Discovery and  
 841 Data Mining, KDD'96. Portland, Oregon: AAAI Press. pp. 226-231.  
 842 Friedman AA, Letai A, Fisher DE, Flaherty KT. 2015. Precision medicine for cancer with next-  
 843 generation functional diagnostics. *Nature Reviews Cancer* **15**:747-756.  
 844 doi:10.1038/nrc4015  
 845 Gilvary C, Madhukar NS, Gayvert K, Foronda M, Perez A, Leslie CS, Dow L, Pandey G, Elemento O.  
 846 2019. A machine learning approach predicts essential genes and pharmacological targets in  
 847 cancer. *bioRxiv* 692277. doi:10.1101/692277  
 848 Haber AL, Biton M, Rogel N, Herbst RH, Shekhar K, Smillie C, Burgin G, Delorey TM, Howitt MR,  
 849 Katz Y, Tirosh I, Beyaz S, Dionne D, Zhang M, Raychowdhury R, Garrett WS, Rozenblatt-  
 850 Rosen O, Shi HN, Yilmaz O, Xavier RJ, Regev A. 2017. A single-cell survey of the small  
 851 intestinal epithelium. *Nature* **551**:333-339. doi:10.1038/nature24489  
 852 Hornik K. 2005. A CLUE for CLUster ensembles. *Journal of Statistical Software* **14**:1-25.  
 853 Ingold I, Berndt C, Schmitt S, Doll S, Poschmann G, Buday K, Roveri A, Peng X, Porto Freitas F, Seibt  
 854 T, Mehr L, Aichler M, Walch A, Lamp D, Jastroch M, Miyamoto S, Wurst W, Ursini F, Arnér  
 855 ESJ, Fradejas-Villar N, Schweizer U, Zischka H, Friedmann Angeli JP, Conrad M. 2018.  
 856 Selenium Utilization by GPX4 Is Required to Prevent Hydroperoxide-Induced Ferroptosis.  
 857 *Cell* **172**:409-422.e21. doi:10.1016/j.cell.2017.11.048  
 858 Jeong W-J, Ro EJ, Choi K-Y. 2018. Interaction between Wnt/ $\beta$ -catenin and RAS-ERK pathways and  
 859 an anti-cancer strategy via degradations of  $\beta$ -catenin and RAS by targeting the Wnt/ $\beta$ -  
 860 catenin pathway. *npj Precision Oncology* **2**:1-10. doi:10.1038/s41698-018-0049-y



861 Jin T, Liu L. 2008. Minireview: The Wnt Signaling Pathway Effector TCF7L2 and Type 2 Diabetes  
862 Mellitus. *Molecular Endocrinology* **22**:2383–2392. doi:10.1210/me.2008-0135

863 Kovačević I, Sakaue T, Majolé J, Pronk MC, Maekawa M, Geerts D, Fernandez-Borja M,  
864 Higashiyama S, Hordijk PL. 2018. The Cullin-3–Rbx1–KCTD10 complex controls endothelial  
865 barrier function via K63 ubiquitination of RhoB. *J Cell Biol* **217**:1015–1032.  
866 doi:10.1083/jcb.201606055

867 Maaten L van der. 2014. Accelerating t-SNE using Tree-Based Algorithms. *Journal of Machine*  
868 *Learning Research* **15**:3221–3245.

869 Maaten L van der, Hinton G. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning*  
870 *Research* **9**:2579–2605.

871 Mass E, Ballesteros I, Farlik M, Halbritter F, Günther P, Crozet L, Jacome-Galarza CE, Händler K,  
872 Klughammer J, Kobayashi Y, Gomez-Perdiguero E, Schultze JL, Beyer M, Bock C, Geissmann  
873 F. 2016. Specification of tissue-resident macrophages during organogenesis. *Science*  
874 **353**:aaf4238. doi:10.1126/science.aaf4238

875 McFarland JM, Ho ZV, Kugener G, Dempster JM, Montgomery PG, Bryan JG, Krill-Burger JM, Green  
876 TM, Vazquez F, Boehm JS, Golub TR, Hahn WC, Root DE, Tsherniak A. 2018. Improved  
877 estimation of cancer dependencies from large-scale RNAi screens using model-based  
878 normalization and data integration. *Nat Commun* **9**. doi:10.1038/s41467-018-06916-5

879 Meyers RM, Bryan JG, McFarland JM, Weir BA, Sizemore AE, Xu H, Dharia NV, Montgomery PG,  
880 Cowley GS, Pantel S, Goodale A, Lee Y, Ali LD, Jiang G, Lubonja R, Harrington WF, Strickland  
881 M, Wu T, Hawes DC, Zhivich VA, Wyatt MR, Kalani Z, Chang JJ, Okamoto M, Stegmaier K,  
882 Golub TR, Boehm JS, Vazquez F, Root DE, Hahn WC, Tsherniak A. 2017. Computational  
883 correction of copy number effect improves specificity of CRISPR–Cas9 essentiality screens  
884 in cancer cells. *Nature Genetics* **49**:1779–1784. doi:10.1038/ng.3984

885 Pan J, Meyers RM, Michel BC, Mashtalir N, Sizemore AE, Wells JN, Cassel SH, Vazquez F, Weir BA,  
886 Hahn WC, Marsh JA, Tsherniak A, Kadoch C. 2018. Interrogation of Mammalian Protein  
887 Complex Structure, Function, and Membership Using Genome-Scale Fitness Screens. *Cell*  
888 *Syst* **6**:555–568.e7. doi:10.1016/j.cels.2018.04.011

889 Sandoval GJ, Pulice JL, Pakula H, Schenone M, Takeda DY, Pop M, Boulay G, Williamson KE, McBride  
890 MJ, Pan J, St. Pierre R, Hartman E, Garraway LA, Carr SA, Rivera MN, Li Z, Ronco L, Hahn WC,  
891 Kadoch C. 2018. Binding of TMPRSS2-ERG to BAF Chromatin Remodeling Complexes  
892 Mediates Prostate Oncogenesis. *Molecular Cell* **71**:554–566.e7.  
893 doi:10.1016/j.molcel.2018.06.040

894 Schenone M, Dančík V, Wagner BK, Clemons PA. 2013. Target identification and mechanism of  
895 action in chemical biology and drug discovery. *Nat Chem Biol* **9**:232–240.  
896 doi:10.1038/nchembio.1199

897 Sergushichev A. 2016. An algorithm for fast preranked gene set enrichment analysis using  
898 cumulative statistic calculation. *bioRxiv* 060012. doi:10.1101/060012

899 Smith I, Greenside PG, Natoli T, Lahr DL, Wadden D, Tirosh I, Narayan R, Root DE, Golub TR,  
900 Subramanian A, Doench JG. 2017. Evaluation of RNAi and CRISPR technologies by large-  
901 scale gene expression profiling in the Connectivity Map. *PLOS Biology* **15**:e2003213.  
902 doi:10.1371/journal.pbio.2003213

903 Sørensen HP, Hedegaard J, Sperling-Petersen HU, Mortensen KK. 2001. Remarkable Conservation  
904 of Translation Initiation Factors: IF1/eIF1A and IF2/eIF5B are Universally Distributed  
905 Phylogenetic Markers. *IUBMB Life* **51**:321–327. doi:10.1080/152165401317190842

- Squires JE, Berry MJ. 2008. Eukaryotic selenoprotein synthesis: Mechanistic insight incorporating new factors and new functions for old factors. *IUBMB Life* **60**:232–235. doi:10.1002/iub.38
- Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP. 2005. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America* **102**:15545–15550. doi:10.1073/pnas.0506580102
- Sulahian R, Kwon JJ, Walsh KH, Pailier E, Bosse TL, Thaker M, Almanza D, Dempster JM, Pan J, Piccioni F, Dumont N, Gonzalez A, Rennhack J, Nabet B, Bachman JA, Goodale A, Lee Y, Bagul M, Liao R, Navarro A, Yuan TL, Ng RWS, Raghavan S, Gray NS, Tsherniak A, Vazquez F, Root DE, Firestone AJ, Settleman J, Hahn WC, Aguirre AJ. 2019. Synthetic Lethal Interaction of SHOC2 Depletion with MEK Inhibition in RAS-Driven Cancers. *Cell Reports* **29**:118–134.e8. doi:10.1016/j.celrep.2019.08.090
- The Cancer Genome Atlas Research Network. 2019. The Cancer Genome Atlas Program. *National Cancer Institute*. <https://www.cancer.gov/tcga>
- Tsherniak A, Vazquez F, Montgomery PG, Weir BA, Kryukov G, Cowley GS, Gill S, Harrington WF, Pantel S, Krill-Burger JM, Meyers RM, Ali L, Goodale A, Lee Y, Jiang G, Hsiao J, Gerath WFJ, Howell S, Merkel E, Ghandi M, Garraway LA, Root DE, Golub TR, Boehm JS, Hahn WC. 2017. Defining a Cancer Dependency Map. *Cell* **170**:564–576.e16. doi:10.1016/j.cell.2017.06.010
- Wang W, Malyutina A, Pessia A, Saarela J, Heckman CA, Tang J. 2019. Combined gene essentiality scoring improves the prediction of cancer dependency maps. *EBioMedicine*. doi:10.1016/j.ebiom.2019.10.051
- Wang X, Wang S, Troisi EC, Howard TP, Haswell JR, Wolf BK, Hawk WH, Ramos P, Oberlick EM, Tzvetkov EP, Ross A, Vazquez F, Hahn WC, Park PJ, Roberts CWM. 2019. BRD9 defines a SWI/SNF sub-complex and constitutes a specific vulnerability in malignant rhabdoid tumors. *Nat Commun* **10**:1–11. doi:10.1038/s41467-019-09891-7
- Wang Y, Probin V, Zhou D. 2006. Cancer therapy-induced residual bone marrow injury- Mechanisms of induction and implication for therapy. *Curr Cancer Ther Rev* **2**:271–279.
- Weiss WA, Taylor SS, Shokat KM. 2007. Recognizing and exploiting differences between RNAi and small-molecule inhibitors. *Nat Chem Biol* **3**:739–744. doi:10.1038/nchembio1207-739

Figure 1

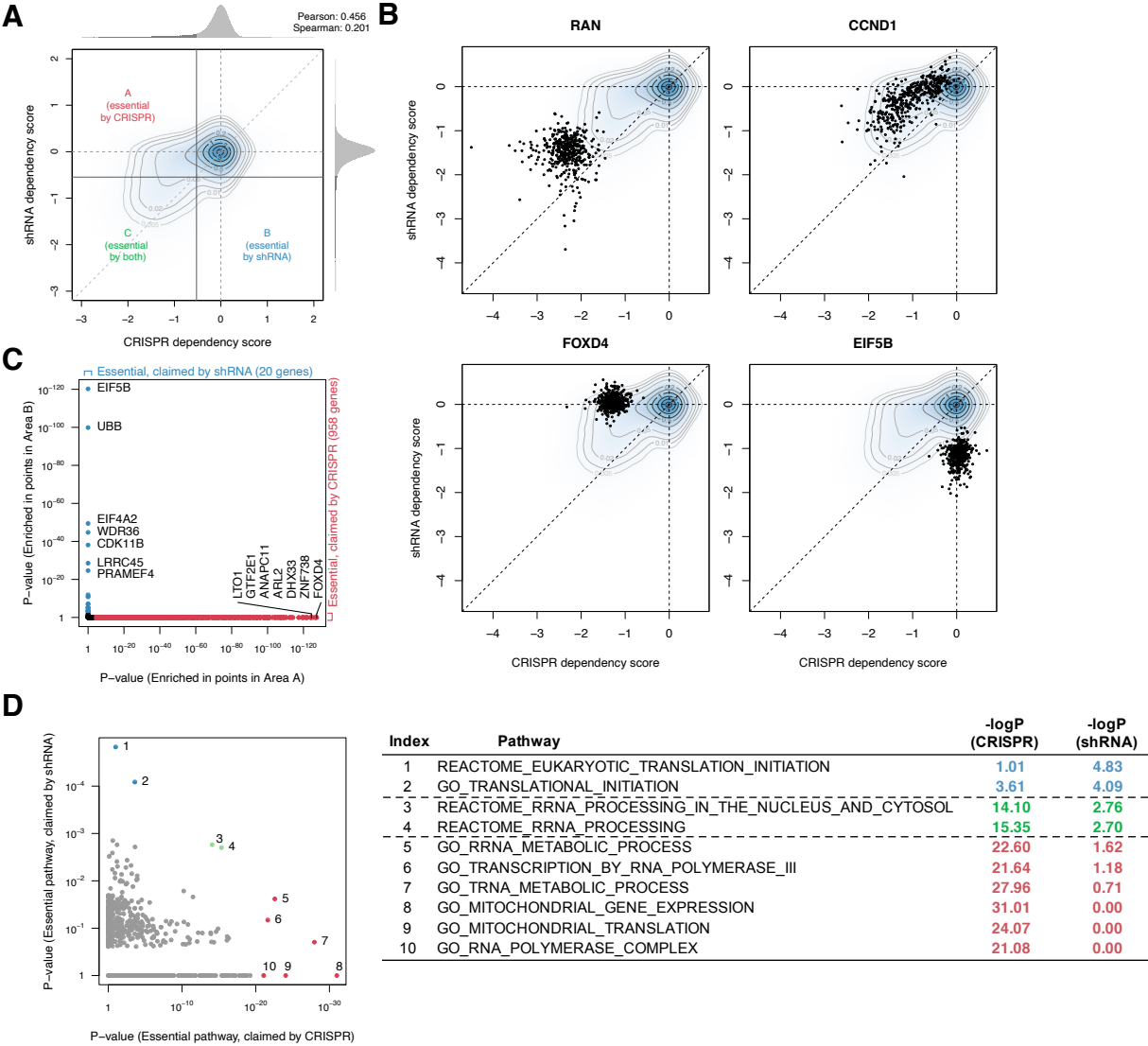


Figure 2

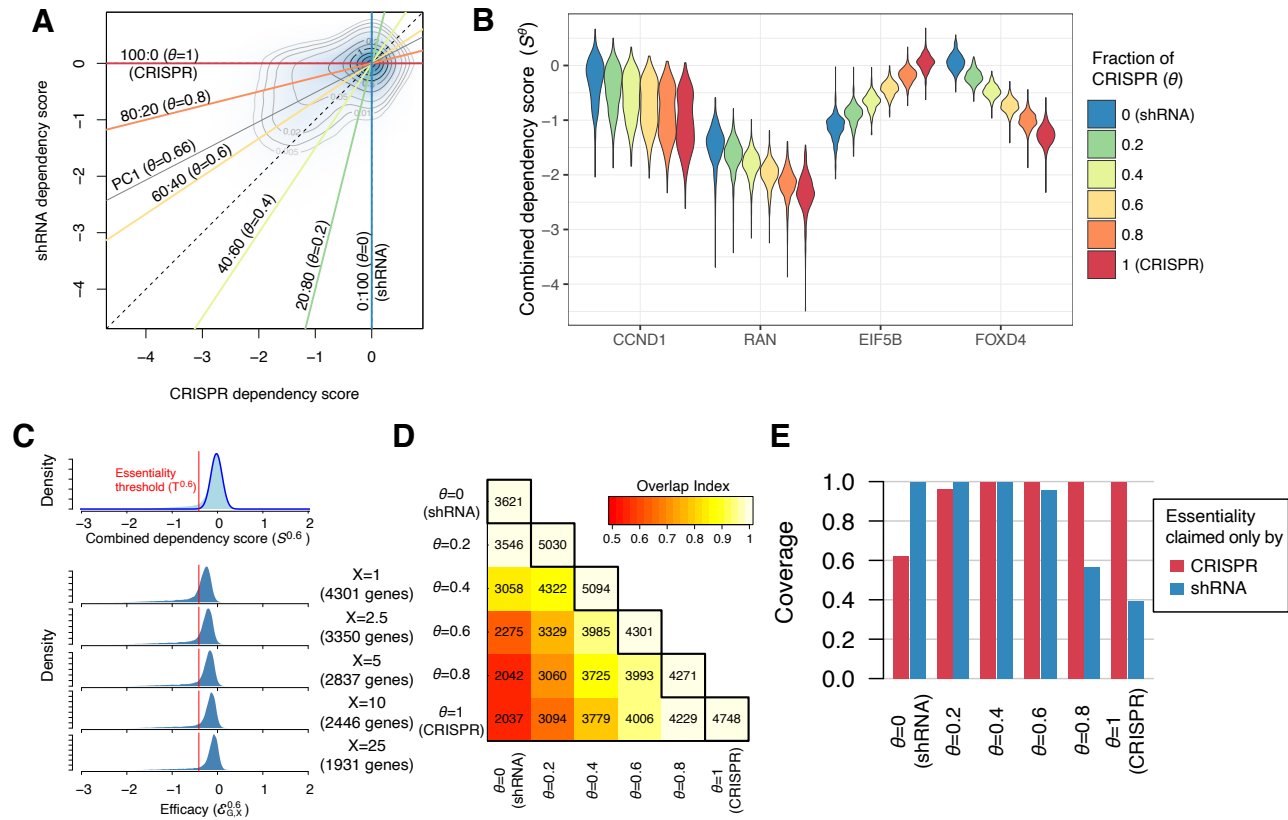


Figure 3

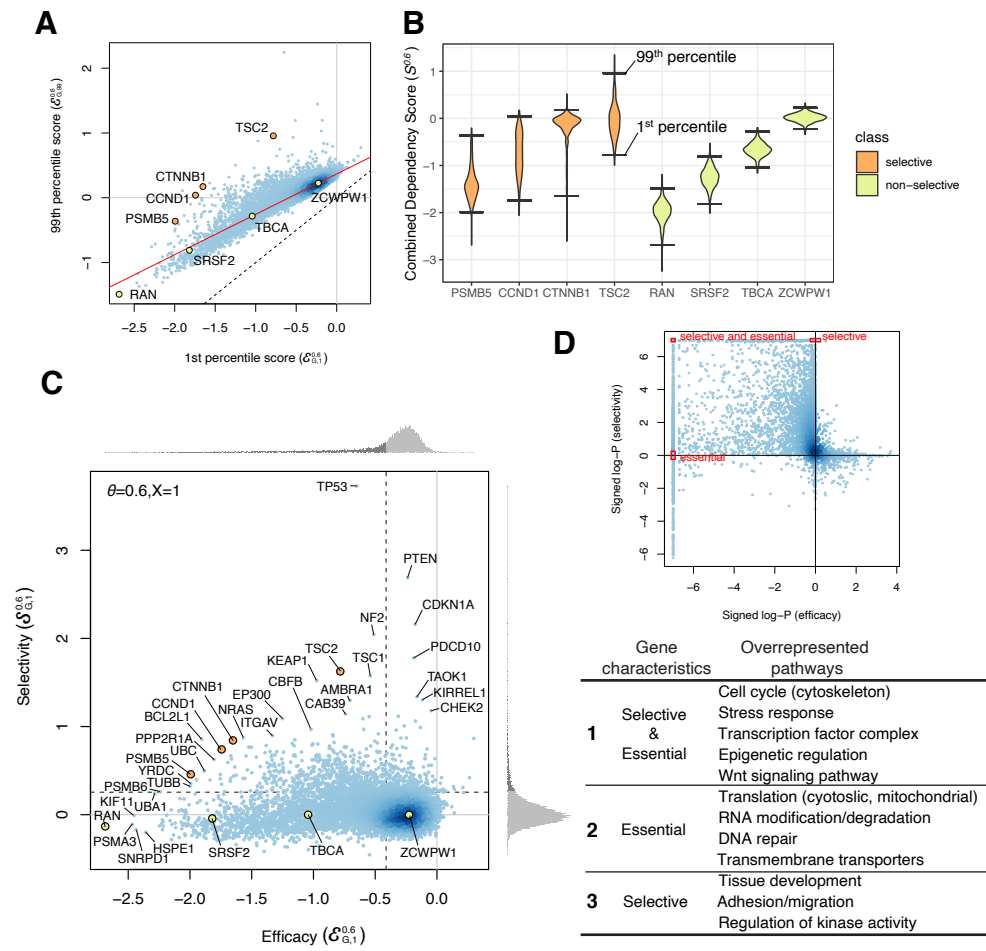


Figure 4

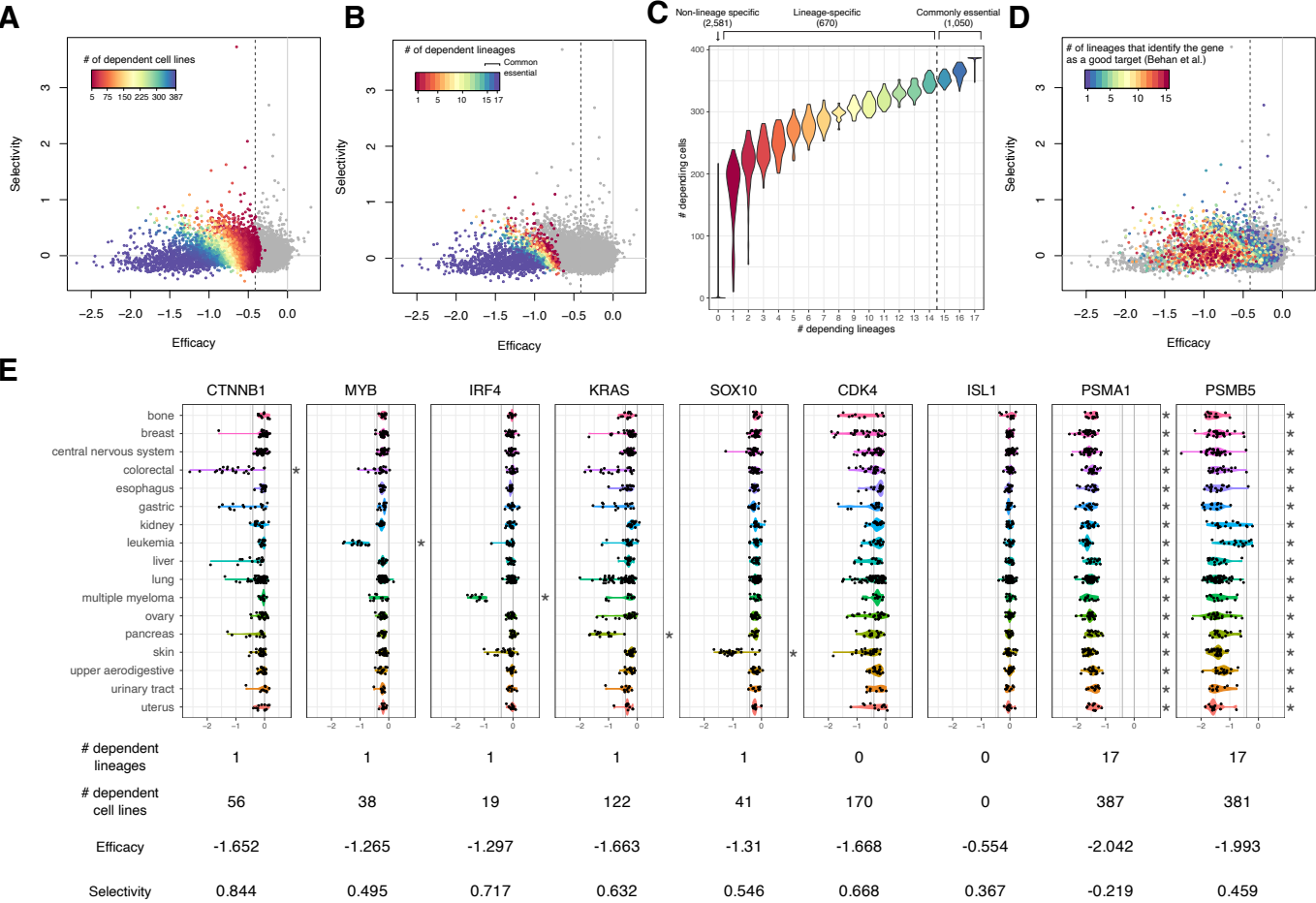


Figure 5

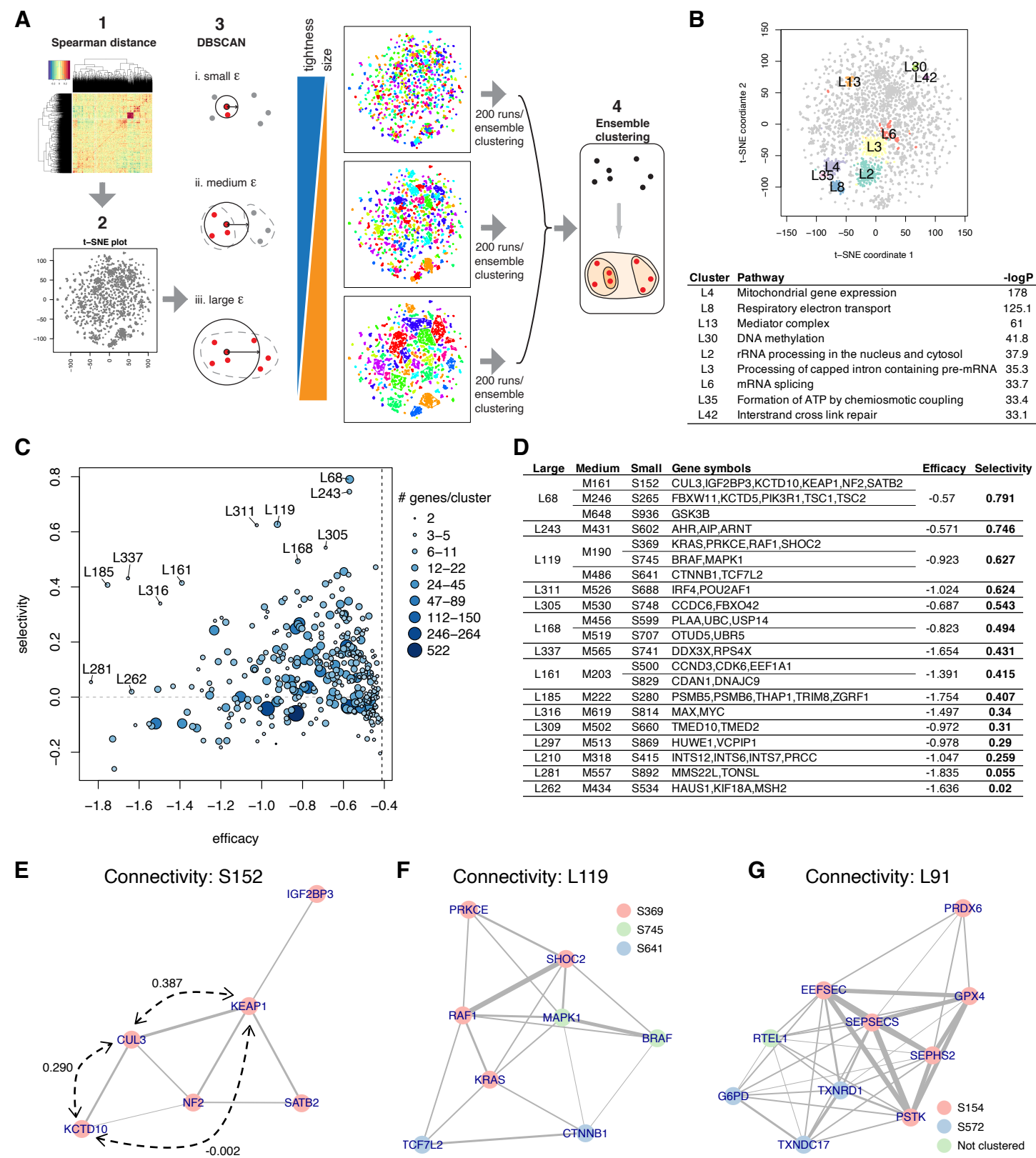
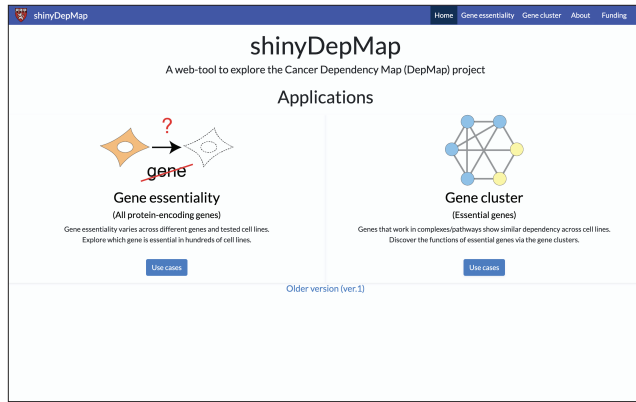
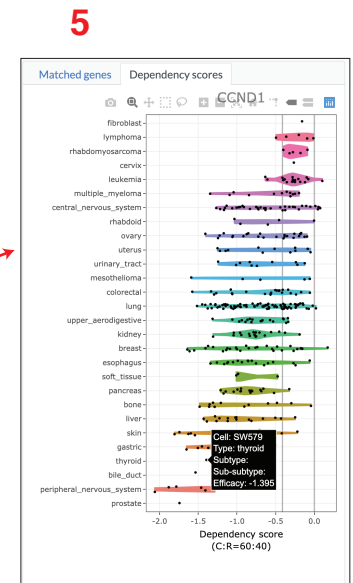
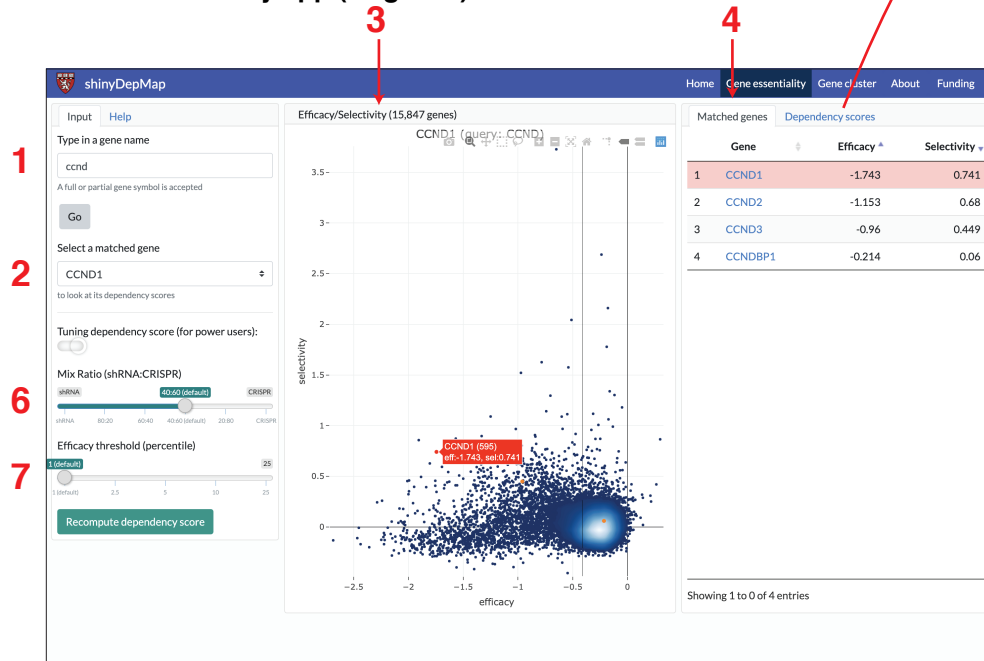


Figure 6

## A Top page



## B. Gene essentiality app (all genes)



## C. Gene clusters app (essential genes only)

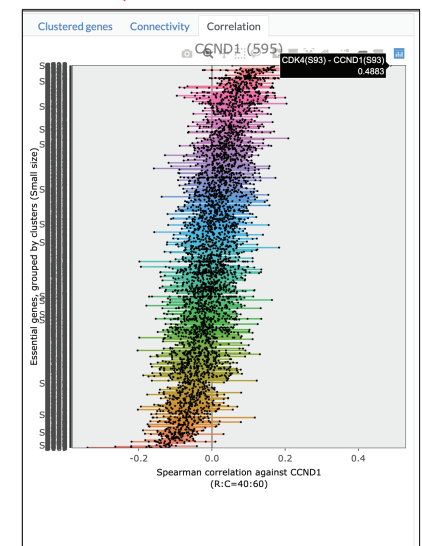
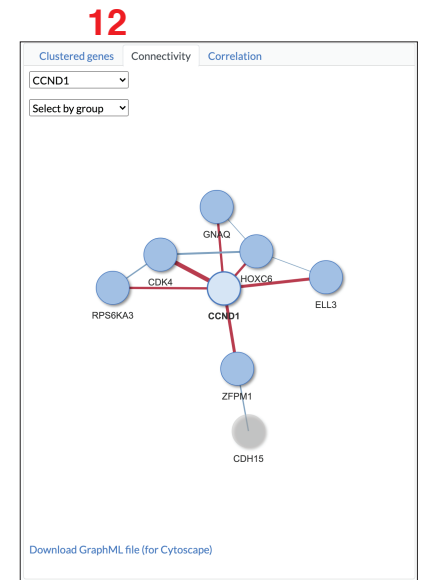
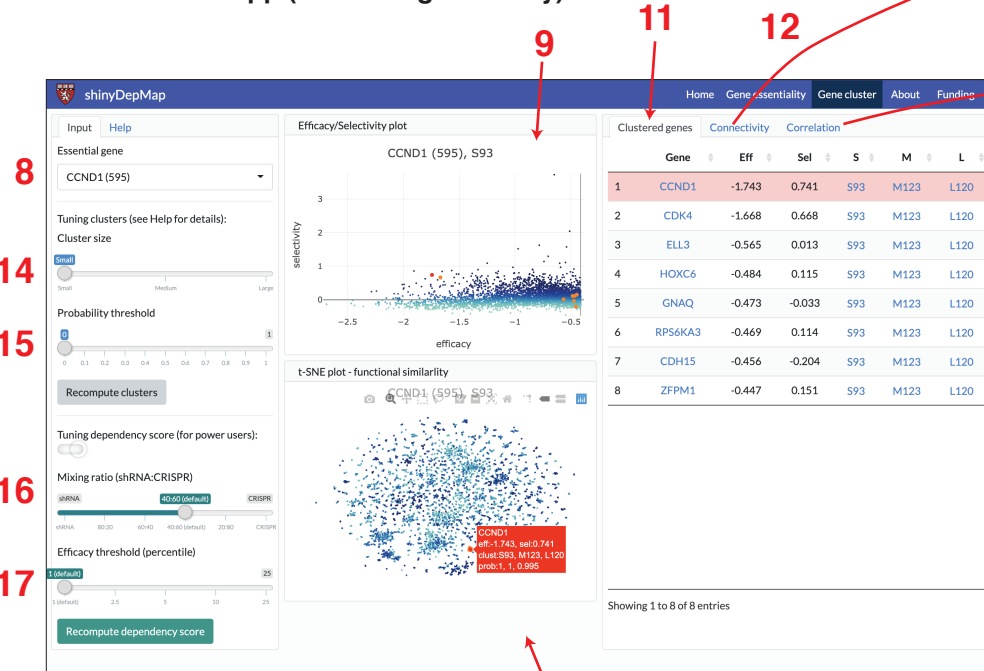




Figure 5 – figure supplement 1

### ## ECHODOTS algorithm

```
Compute the Spearman distance  $\mathbf{D}$  across all the essential genes ----- (1)
For each  $\mathbf{d}$  in 30:200
  For each  $\mathbf{i}$  in 1:200
    Set  $\mathbf{i}$  as the initial seed
    Run t-SNE with  $\mathbf{D}$  as density matrix and get 2D coordinates ----- (2)
    Compute the range of the data points in t-SNE map  $\mathbf{L}(\mathbf{i})$ 
    Compute the neighborhood threshold  $\varepsilon(\mathbf{i}) = \mathbf{L}(\mathbf{i})/\mathbf{d}$  ----- (3)
    Run DBSCAN with  $\varepsilon(\mathbf{i})$  to find clusters for the essential genes  $\mathbf{C}(\mathbf{d}, \mathbf{i})$  ----- (4)
    Find consistent clusters  $\mathbf{CC}(\mathbf{d})$  among  $\mathbf{C}(\mathbf{d}, \mathbf{i})$ ,  $\mathbf{i} = \{1, \dots, 200\}$  ----- (5)
Find the lower bound  $\mathbf{d}_0$  such that  $\mathbf{d} \geq \mathbf{d}_0$  ----- (6)
Pick three  $\mathbf{d}$  and corresponding cluster sets  $\mathbf{CC}(\mathbf{d})$  ----- (7)
```

Figure 5 – figure supplement 2

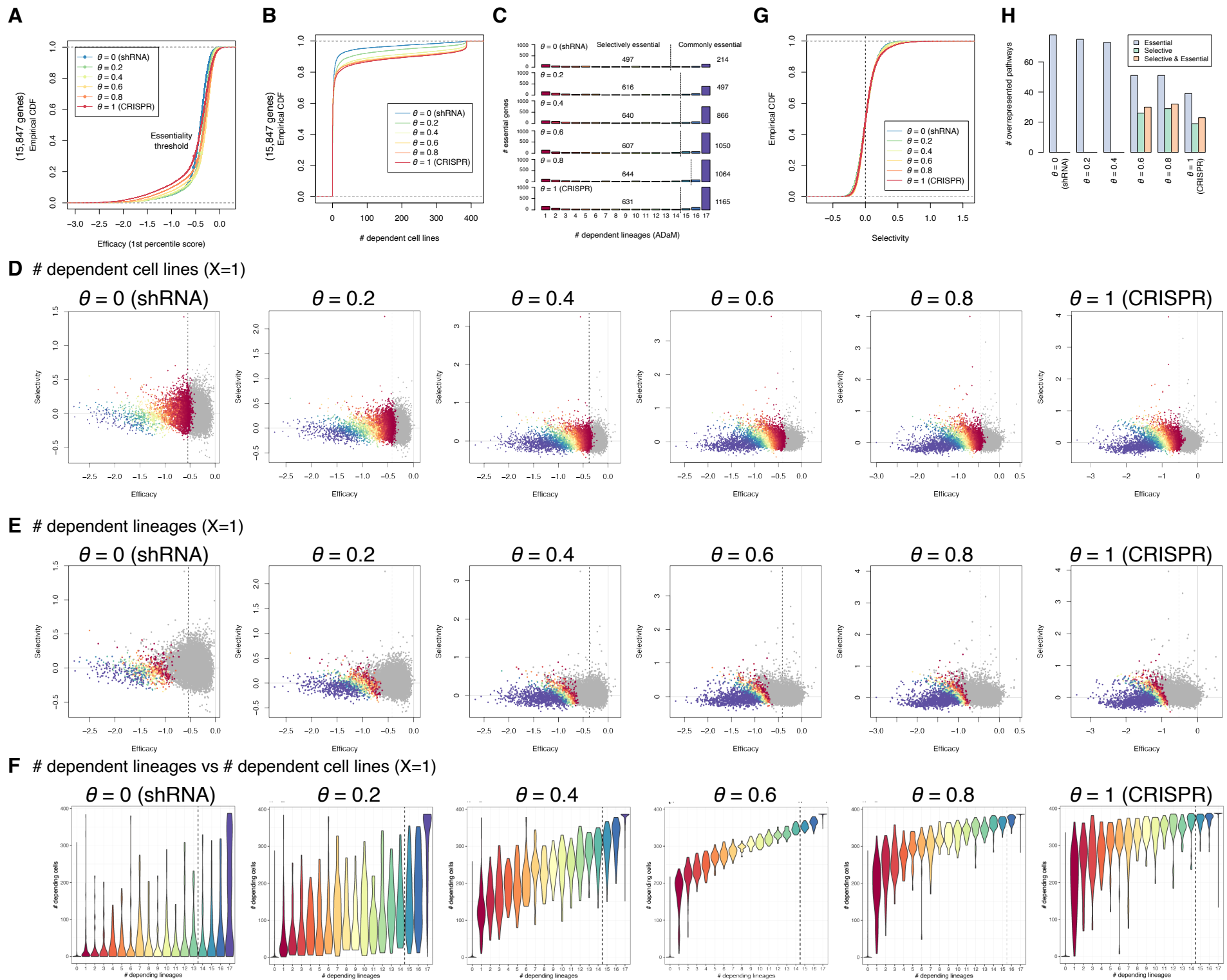


Figure 5 - figure supplement 3

