



Beta-binomial modeling of CRISPR pooled screen data identifies target genes with greater sensitivity and fewer false negatives

Hyun-Hwan Jeong, Seon Young Kim, Maxime W.C. Rousseaux, et al.

Genome Res. published online April 23, 2019

Access the most recent version at doi:[10.1101/gr.245571.118](https://doi.org/10.1101/gr.245571.118)

P<P	Published online April 23, 2019 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Beta-binomial modeling of CRISPR pooled screen data** 2 **identifies target genes with greater sensitivity and fewer false** 3 **negatives**

4 Hyun-Hwan Jeong^{1,2}, Seon Young Kim^{1,2}, Maxime W. C. Rousseaux^{1,2,+}, Huda Y. Zoghbi^{1,2,3,4},
5 Zhandong Liu^{2,3*}

6 ¹Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas, USA

7 ²Jan and Dan Duncan Neurological Research Institute, Texas Children's Hospital, Houston, Texas, USA

8 ³Department of Pediatrics, Baylor College of Medicine, Houston, Texas, USA

9 ⁴Howard Hughes Medical Institute, Houston, Texas, USA

10 ⁺Present Address: University of Ottawa Brain and Mind Research Institute and Department of Cellular
11 and Molecular Medicine, University of Ottawa, Ottawa, Ontario, Canada

12 ^{*}Correspondence should be addressed to: zhandong.liu@bcm.edu

13 **Abstract**

14 The simplicity and cost-effectiveness of CRISPR technology have made high-throughput pooled
15 screening approaches accessible to virtually any lab. Analyzing the large sequencing data derived from
16 these studies, however, still demands considerable bioinformatics expertise. Various methods have been
17 developed to lessen this requirement, but there are still three tasks for accurate CRISPR screen analysis
18 that involve bioinformatic know-how if not prowess: designing a proper statistical hypothesis test for
19 robust target identification, developing an accurate mapping algorithm to quantify sgRNA levels, and
20 minimizing the parameters necessary that need to be fine-tuned. To make CRISPR screen analysis more
21 reliable as well as more readily accessible, we have developed a new algorithm, called
22 CRISPRBetaBinomial or CB² (<https://CRAN.R-project.org/package=CB2>). Based on the beta-binomial
23 distribution, which is better suited to sgRNA data, CB² outperforms the eight most commonly used
24 methods (HiTSelect, MAGeCK, PBNPA, PinAPL-Py, RIGER, RSA, ScreenBEAM, and sgRSEA) in
25 both accurately quantifying sgRNAs and identifying target genes, with greater sensitivity and a much
26 lower false discovery rate. It also accommodates staggered sgRNA sequences. In conjunction with
27 CRISPRcloud, CB² will bring CRISPR screen analysis within reach for a wider community of
28 researchers.

29 **Introduction**

30 Genetic screens have become a favored tool for gathering information about disease pathogenesis and
31 cellular biology. Initially, these screens were performed using chemical mutagenesis or RNA interference
32 (RNAi), which are effective but laborious processes (Mohr et al. 2014; Park et al. 2013; Schlabach et al.
33 2008; Silva et al. 2008; Simon et al. 2015). Larger-scale, pooled approaches were finally made feasible by
34 the advent of microarray (DeJesus et al. 2016; Gilbert et al. 2014; Luo et al. 2009; Shalem et al. 2014;
35 Wang et al. 2014). Pooled shRNA (short-hairpin RNA) libraries can be barcoded and packaged into
36 viruses, which are used to infect a population of cells that are then selected for a desired phenotype (e.g.,
37 growth or fluorescence). Hybridizing microarray soon followed for hit identification (Paddison et al.
38 2004). In later iterations of this approach, next-generation sequencing (NGS) was used to identify hits
39 (Hu and Luo 2012).

40 The development and optimization of CRISPR/Cas9 (Clustered Regularly Interspaced Short Palindromic
41 Repeats and CRISPR-associated protein 9) systems have propelled pooled screen approaches into even
42 wider use. Besides their relatively simplicity and low cost, the robustness of hit identification has reduced
43 the requirement for redundancy in the number of targeting single-guide RNAs (sgRNAs), which allows
44 the same size library to be more diverse (Sanjana et al. 2014; Xu et al. 2015). Moreover, these pooled
45 libraries have been made accessible through repositories such as Addgene (<https://www.addgene.org/>).
46 CRISPR/Cas9 pooled screens are thus within the technical reach of most biomedical researchers (Doench
47 2017). For all this accessibility on the experimental side, however, the resulting bioinformatics datasets
48 are enormous and complex. The analysis of thousands of genetic perturbations demands considerable
49 bioinformatics expertise.

50 In our experience, most users are unaware of whether their tool of choice models the data according to
51 Poisson, negative binomial, or even Gaussian distribution, nor what the relative strengths and weaknesses
52 of these models are. The current roster of tools bear the imprint of the history of RNA-seq data analysis:
53 RNA-seq data were initially modeled using Poisson distributions (Marioni et al. 2008), which is a natural
54 choice for simple read counts. Poisson distribution assumes that the mean and variance are equal,
55 however, and with biological data, the variance is often greater than the mean. Analytic methods therefore
56 turned to negative binomial distribution, which can handle overdispersed data (Anders and Huber 2010;
57 Love et al. 2014). A number of popular tools still use negative binomial distribution to analyze sgRNA
58 screen data (Li et al. 2014; Spahn et al. 2017), even though the structure of the sgRNA screen data is very
59 different from that of RNA-seq. In the latter, there is huge variation in transcript lengths, from 60bp to
60 2.4Mbp, but all sgRNAs for any given gene are designed to have the same length. This often leads to the
61 variance being less than the mean (Supplemental Figure 1). We hypothesized that a beta-binomial model,

62 in which the variance can be either larger or smaller than the mean, would better fit the data and more
63 accurately identify changes in sgRNA.

64 We therefore developed CB², a new web-based algorithm that uses beta-binomial distribution, and
65 compared its performance with that of the eight most commonly used algorithms (HiTSelect (Diaz et al.
66 2015), MAGeCK (Li et al. 2014), PBNPA (Jia et al. 2017), PinAPL-Py (Spahn et al. 2017), RIGER (Luo
67 et al. 2009), RSA (König et al. 2007), ScreenBEAM (Yu et al. 2016), and sgRSEA (Noh 2015)), which
68 encompass both parametric and nonparametric approaches (Table 1). We applied all these methods to ten
69 different biological datasets, taken from fields ranging from cancer to basic cell biology (Koike-Yusa et
70 al. 2014; Parnas et al. 2015; Evers et al. 2016; Golden et al. 2017; Li et al. 2018; Sanson et al. 2018).

71 Results

72 CB² is more sensitive in target gene identification than existing methods

73 Identifying candidates by a statistical hypothesis test is a key component in any screen analysis. In CB²,
74 we adapted a beta-binomial model (Baggerly et al. 2003) with a modified Student's *t*-test to measure
75 differences in single-guide RNA (sgRNA) levels, followed by Fisher's combined probability test (Fisher
76 1925) to estimate the gene-level significance. We chose Fisher's method for two reasons: first, to keep the
77 entire pipeline parametric, and second, to keep CB² as fast as possible (Robust Rank Aggregation [RRA]
78 requires permuting the data—a non-parametric feature—so it runs slower than Fisher's method).
79 Furthermore, when we compared Fisher's method against RRA, we found that RRA is not effective in
80 combining the p-values estimated by beta-binomial distribution (Supplemental Figure 2).

81 We compared CB² with eight state-of-the-art methods on three benchmark datasets evaluating gene
82 essentiality (Evers et al. 2016) using different technologies: CRISPRn (CRISPR nuclease gene knockout
83 via Cas9) and CRISPRi (CRISPRinterference, a CRISPR/Cas9 system with a catalytically inactive Cas9
84 fused to the transcriptional repressor KRAB which results in gene repression). These benchmark datasets
85 (CRISPRn-RT112, CRISPRn-UMUC3, and CRISPRi-RT112) were constructed based on 46 genes that
86 are essential for cell survival and 47 genes that are non-essential. We first tested whether these methods
87 could easily distinguish between essential and non-essential genes. We found that each method clearly
88 discriminates essentiality by their gene rankings (Supplemental Figure 3A). In addition, gene rankings
89 obtained from each method, except PBNPA for the CRISPRi-RT112 dataset, are highly correlated
90 (Supplemental Figure 3B , R^2 is [0.86,0.98] for CRISPRn-RT112, [0.85,0.98] for CRISPRn-UMUC3,
91 and [0.72,0.96] for CRISPRi-RT112). CB², ScreenBEAM, and MAGeCK produced very similar gene
92 rankings across all the benchmark datasets. We also compared the Precision-Recall (PR) and Receiver

93 operating characteristic (ROC) curves across the methods and calculated the Area under the curve (AUC)
94 of each. CB² recorded the best PR-AUC and ROC-AUC scores for both CRISPRn screen datasets
95 (Supplemental Figure 4 & 5). HiTSelect had the best PR-AUC and ROC-AUC scores for the CRISPRi-
96 RT112 dataset for which CB² achieved comparable scores (Supplemental Figure 6). Although the gene
97 ranking is similar among these methods (Supplemental Figure 3), the estimated p-values and FDRs are
98 very different. These results highlight the importance of using FDR to guide the gene selection process.

99 While several CRISPR screens (Aguirre et al. 2016; Parnas et al. 2015; Zhou et al. 2014) have prioritized
100 candidate hits by ranking, they do not provide statistical estimates of error rates. These methods,
101 therefore, rely on an arbitrary rule to select the top candidate genes and are prone to biased selections and
102 high hit attrition rates. One solution is candidate selection by a quantitative statistical measure such as a p-
103 value or a false positive rate (FDR) cut-off. To assess the detection powers of FDR of established
104 CRISPR screen analysis methods and CB², we measured the F1-score ($2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$) of each
105 method, i.e., the harmonic average of the precision and the recall, for FDR thresholds ranging from 10%
106 to 0.01%. CB² outperformed all other methods at every FDR cut-off level, and all other methods lost their
107 detection powers at more rigorous FDRs (Figure 1 A&B and Supplemental Figure 7). In other words, all
108 methods demonstrated a small type-I error due to the strong lethality phenotype of the CRISPR assay, but
109 CB² demonstrated a significantly lower type-II error than the other methods (Supplemental Figure 8).
110 Across all paradigms tested with different FDR cut-offs, CB² performed the best, with a much larger F1-
111 score and recall. Thus, CB² is both sensitive and specific in selecting candidate genes.

112 To understand the differences produced by these methods, we next tested a prototypical essential gene,
113 *RPL5*, to compare the gene-level enrichment across datasets and analytical tools. In the first CRISPR
114 screen on an RT112 cell line, we expected to see the depletion of sgRNAs targeting *RPL5* in group T_1 .
115 Out of the 10 sgRNAs that target this gene in the CRISPRn-RT112 dataset, six showed a strong decrease
116 in the group T_1 . CB² estimated an FDR of 2.07×10^{-19} , while only three other methods (HiTSelect,
117 ScreenBEAM, and sgRSEA) estimated FDR to be < 0.01 (Figure 2A). Next, we looked at the same gene
118 in the UMUC3 cell line. Five of the 10 sgRNAs targeting *RPL5* decreased in group T_1 , and all five
119 sgRNAs were among those identified in RT112 cell line. CB² estimated an FDR of 3.81×10^{-10} for
120 *RPL5*, while none of the other methods identified it to be significantly depleted with an FDR cutoff below
121 0.01 (Figure 2B). Lastly, in the CRISPRi-RT112 dataset, three of seven sgRNAs indicated depletions, but
122 only CB² was able to estimate the FDR of 2.78×10^{-8} , and the other methods did not count *RPL5* as a
123 hit in the dataset (Figure 2C). Overall, CB² produced more reliable hit identification than other methods
124 based on statistical cut-offs for the gene tested (*COPS8* and *RPL27*, Supplemental Figure 12 & 13).

125 We performed the same analysis on two distinct datasets (Sanson et al. (2018)) to determine how CB^2
126 performs compared to other methods for genome-wide screening analysis. Sanson et al. used novel
127 optimized libraries for genome-wide CRISPRn (Brunello), CRISPRi (Dolcetto), and CRISPRa
128 (Calabrese) screening and showed these libraries outperform other previously established libraries, such
129 as GeCKO (Sanjana et al. 2014) and hCRISPRi-v2 (Horlbeck et al. 2016). For our analysis, we chose two
130 screens for benchmarking, both datasets from a dropout screen in A375 melanoma cells: one used the
131 Brunello CRISPRn library with tracr-v2 tracrRNA (CRISPRn-A375), the other used the Dolcetto
132 CRISPRi Set A library (CRISPRi-A375). Each dataset contains a control sample (plasmid DNA) and
133 three biological replicates. We used the gold-standard gene sets of 1,580 essential and 927 non-essential
134 genes reported in (Hart et al. 2014, 2015) to assess the performance of the methods. (We excluded
135 PinAPL-Py from benchmarking since it does not report statistics when the input contains only one control
136 sample.) CB^2 outperformed other methods at the stringent FDR cutoff level (Figure 2 and Supplemental
137 Figure 9). CB^2 outperformed all other methods in F-1, precision and recall measures at the stringent FDR
138 cut-offs on (Sanson et al. 2018)'s A375 genome-wide screen datasets. F1-score (top), precision (middle),
139 and recall (bottom) for each method on two benchmark datasets are presented as a function of FDR cut-
140 off values (Figure 1 C & D and Supplemental Figure 9), and provided higher AUC values of PR and ROC
141 curves than the other methods (Supplemental Figure 10 & 11).

142 These results indicate that CB^2 more accurately estimates the gene-level FDR. The use of FDR in
143 selecting hits is critical in real data analysis since the arbitrary selection of top genes is purely heuristic.
144 CB^2 is better at identifying true hits based on multiple concordant sgRNAs targeting the same gene.

145 **CB^2 is more sensitive in target gene detection than existing methods**

146 To test the idea that a beta-binomial model would better fit the data and more accurately identify changes
147 in sgRNA, we compared the sgRNA level statistics on several CRISPR pooled libraries containing non-
148 targeting sgRNAs as negative controls. Non-targeting sgRNAs are not supposed to show any differential
149 enrichment and can be used to assess the quality of the method. (Parnas et al. 2015) used Mouse CRISPR
150 Knockout Pooled Library for their genome-wide screen (GeCKO v2, Addgene #1000000052,
151 #1000000053), which contains 1,000 non-targeting sgRNAs. We therefore used this dataset to measure
152 the specificity with which CB^2 and MAGeCK detect true negatives. We compared the unadjusted p-
153 values for sgRNAs since the FDR is controlled at the gene level.

154 CB^2 showed greater specificity (the proportion of actual negatives that are correctly identified) than
155 MAGeCK across a wide range of p-value thresholds. At a p-value threshold of 0.01, CB^2 had a specificity
156 of 86% vs. MAGeCK's 68% (Figure 3A). Next, we plotted the log-fold change against the p-value levels

157 in a volcano plot. The majority of the negative control sgRNAs were correctly detected by CB², whereas
158 MAGeCK demonstrated a one-side long tail for positively changed sgRNAs producing inflated p-values
159 for a group of negative controls (Figure 3B). Many of these false positives showed extremely low p-
160 values (ranging from 10⁻⁵ to 10⁻⁴⁰), indicating a strong selection bias. To understand the impact of this
161 selection bias, we analyzed the rest of the sgRNA library with both methods. At the same threshold
162 ($p < 0.01$), CB² selected 12,648 sgRNAs while MAGeCK selected 31,381 sgRNAs; 2,971 sgRNAs were
163 identified by both methods (Figure 4B). We applied the same analysis to the CRISPRn-A375 dataset
164 (Sanson et al. 2018) which contains 1,000 non-targeting sgRNAs. CB² shows higher specificity than
165 MAGeCK, except when setting a p-value cut-off at 0.2. Similar p-value distributions shown in Figure 3
166 were also found for this dataset (Supplemental Figure 14).

167 We plotted sgRNAs unique to each method on a heatmap, which showed a high concordance within each
168 experimental group for sgRNAs unique to CB² (Figure 4A). In contrast, sgRNAs identified by MAGeCK
169 showed a much noisier pattern, and samples from the same experimental group could not be clustered
170 together based on these differentially enriched sgRNAs (Figure 4A).

171 Next, we performed the same sgRNA-level comparisons on two additional datasets. In the first study, a
172 differentiation screen was conducted to identify target genes that maintain naive pluripotency (Li et al.
173 2018) using the Mouse Improved Genome-wide Knockout CRISPR Library v2 (Addgene #67988). The
174 library contains 91,319 sgRNAs targeting 18,542 mouse genes. To identify differentially enriched
175 sgRNAs, we kept the same threshold ($p < 0.01$) for both CB² and MAGeCK. CB² identified 732 sgRNAs
176 while MAGeCK identified 5,105 sgRNAs (395 sgRNAs shared between the two) (Figure 4C and 4D),
177 and we observed the same trend as in Parnas et al.'s screening dataset (Parnas et al. 2015) (Figure 4A and
178 B). We found the same trend on the Evers et al.'s datasets (Evers et al. 2016) (Supplemental Figure 15).
179 Thus, CB²'s accurate sgRNA-level statistics are attributable to its use of the beta-binomial model.

180

181 **CB² provides more accurate alignment without parameter tuning**

182 Many CRISPR pooled screens use in-house scripts to quantify sgRNA abundance (Gilbert et al. 2014;
183 Golden et al. 2017; Iorio et al. 2018; Li et al. 2018) or other alignment algorithms for RNA-seq (DeJesus
184 et al. 2016; Hart et al. 2015; Parnas et al. 2015; Sanjana et al. 2014). These codes are often not shared
185 publicly and are not easily reusable. Both MAGeCK and PinAPL-Py provide an integrated mapping
186 function, but PinAPL-Py requires complex parameter tuning and MAGeCK samples only the first million
187 reads to estimate the location of sgRNAs in FASTQ files. Furthermore, there is no systematic comparison
188 of mapping accuracy in the literature, and users lack reliable guidelines for selecting mapping tools. We,

189 therefore, introduced an adaptive hash-mapping algorithm into CB² and tested all three methods on six
190 published datasets (Supplemental Table 1).

191 CB² showed consistently greater mapping accuracy than MAGeCK and PinAPL-Py (Figure 5A). To
192 understand why, we studied the reads that are mapped by CB² but not MAGeCK or PinAPL-Py in the
193 CRISPRn-RT112 dataset (Evers et al. 2016). MAGeCK mapped 64% of the reads, compared to 75% by
194 CB² (Figure 5A). This is primarily due to the fact that MAGeCK estimates the trimming windows using
195 the first *N* reads from the input (*N* is 100,000 by default). There is no guarantee that these windows are
196 optimal for the rest of the input files. If a sgRNA locates outside of the precomputed windows, MAGeCK
197 will fail to detect it (Figure 5B). PinAPL-Py does not precompute sgRNA locations based on a subset of
198 reads but uses Cutadapt (Martin 2011) for flexible trimming followed by the Bowtie 2-based alignment
199 (Langmead and Salzberg 2012). We found that PinAPL-Py failed to identify some of the sgRNAs because
200 reads fail to align due to the incorrect trimming from Cutadapt even under several different tuning
201 parameters (Figure 5B). This is likely because of the frequent indels that occurred in the 5' adapter
202 sequence region of the reads (see Figure 5B)—usually, the quad-nucleotide sequence ‘CACC’, which is
203 part of the U6 promoter, precedes the sgRNA sequence. In contrast, the ‘GTTT’ sequence, which is the
204 first four nucleotides of the sgRNA scaffold sequence, was present in all the reads. Given the fidelity of
205 the GTTT sequence, the sequences mapped by CB² but missed by other algorithms are likely accurate
206 and not false positives. CB² is currently the only CRISPR/Cas9 online screen analysis tool with
207 parameter-free mapping and high accuracy in sgRNA quantification.

208

209 **CB² is accessible, secure, and does not require a steep learning curve with** 210 **CRISPRcloud**

211 We had previously developed a web-based application called CRISPRcloud that could run any statistical
212 testing and mapping algorithm through the cloud-based infrastructure provided by Amazon Web Service
213 (Jeong et al. 2017). We implemented CB² in the platform and added new features to increase speed and
214 data security (Supplemental Figure 16). CRISPRcloud is compatible with most modern web browsers
215 (Google Chrome version 69 and later, Apple Safari version 11 and later, and Mozilla Firefox 48.0 and
216 later) and operating systems (iOS, Windows, and Linux). Our fast client-side sgRNA mapping program
217 reduces input files of several gigabytes into a single megabyte-sized file. By transferring a much smaller
218 file through the Internet, CRISPRcloud decreases transfer time and prevents the sharing of raw input files,
219 thereby eliminating downloading errors and data privacy issues in one step. Our adaptive mapping
220 algorithm, via Angular (<https://angular.io/>) and TypeScript (<https://www.typescriptlang.org>), provides an
221 open-source front-end web application platform. The enormous computing power needed to perform

222 these operations mean that platforms built with a centralized server solution will have load-balancing
223 problems when many users submit their requests simultaneously, leading to the longer user waiting times
224 and raising the risk of system-wide failure. CB² therefore provides a decentralized, cloud-computing-
225 based, scalable service through a combination of AWS infrastructure that includes Amazon Elastic
226 Compute Cloud (EC2) (<https://aws.amazon.com/ec2/>), Amazon Simple Storage Service (S3)
227 (<https://aws.amazon.com/s3/>) and Amazon Simple Queue Service (SQS) (<https://aws.amazon.com/sqs/>).
228 With this infrastructure, we launched a web service of CRISPRcloud. CRISPRcloud enables researchers
229 with no programming background to pre-process, check the quality, statistically analyze, query, and
230 visualize their CRISPR/Cas9 pooled screening data (Supplemental Table 2).

231 **Discussion**

232 The number of datasets for CRISPR/Cas9 screens in Gene Expression Omnibus have more than tripled
233 each of the past three years (39 datasets in 2015, 121 datasets in 2016, and 408 datasets in 2017). This
234 expansion has outpaced the development of methods for analyzing the data, most of which employ
235 statistical models that are better suited to RNA-seq than to sgRNA data. Here we took into account the
236 difference between the two types of data to develop a new algorithm, CB², and show that it is more
237 sensitive, specific, and selective than eight other leading tools.

238 We focused first on the central task for any analytic tool being applied to sgRNA data: statistical
239 hypothesis testing to identify target genes accurately from the screening data. Several methods that
240 facilitate analysis of RNAi pooled screening data (Dutta et al. 2016; König et al. 2007; Luo et al. 2008;
241 Shao et al. 2013) are not compatible with CRISPR/Cas9 pooled screening data because of differences in
242 effect size, sequence determinants, and on- vs. off-target effects (Li et al. 2014). MAGeCK was the first
243 tool specifically developed to analyze CRISPR/Cas9 pooled screening data, and it combines a negative-
244 binomial distribution model with a modified robust ranking aggregation (RRA) algorithm (Li et al. 2014).
245 Subsequent methods (Table 1) employed different strategies to improve the accuracy of data analysis, but
246 to our knowledge there has never been a thoroughgoing attempt to benchmark these methods and
247 determine which performs best with sgRNA data. Our choice of the beta-binomial distribution, which is
248 not the approach used by any of these analytic tools, was justified by both the theoretical and empirical
249 considerations and proved able to provide far fewer false positives than these other methods at
250 comparable FDR thresholds.

251 CB² also addressed another difficult task: quantifying sgRNA from next-generation sequencing data.
252 Except for the quantification algorithm provided by MAGeCK, most studies use in-house algorithms or

253 extend established methods that were optimized for RNA-seq. CB² proved capable of fast and accurate
 254 alignment, with the ability to handle indels.

255 Last but not least is the challenge of making powerful tools readily accessible to the research community.
 256 Of the existing tools, PinAPL-Py (Spahn et al. 2017) and CRISPRcloud (Jeong et al. 2017) are the only
 257 two that support a graphical web interface and require no additional program installation. These programs
 258 are an important first step toward enabling the scientists who are actually generating the CRISPR/Cas9
 259 screen data to analyze their large datasets. They still have limitations, however: for PinAPL-Py, users still
 260 need to provide the adaptor sequence to be trimmed, the error tolerance rate, and the quality threshold for
 261 trimmed reads. CRISPRcloud is the only framework into which the user can plug in any statistical tools
 262 or mapping algorithms but transferring a large amount of sequence data over the internet has inherent
 263 disadvantages such as long transfer times, vulnerability to file copying errors, and possible data security
 264 breaches. Taking advantage of the CRISPRcloud framework, CB² is fully web-based and designed to
 265 require only the minimal number of parameters for data analysis, since fewer parameters means shorter
 266 learning curves for the majority of users. CB²'s power and accessibility will enable more labs to extract
 267 biologically relevant discoveries from CRISPR pooled screens.

268 **Methods and Materials**

269 **Statistical hypothesis testing using beta-binomial distribution for sgRNA-level** 270 **differential analysis**

271 We adapted a beta-binomial model proposed for Serial Analysis of Gene Expression (SAGE) by
 272 (Baggerly et al. 2003). Specifically, let p_i be the true proportion of an sgRNA in sample i . We assume the
 273 value of p_i can vary from sample to sample and follows a beta distribution, $p_i \sim Beta(\alpha, \beta)$. Let X_i
 274 denote the number of read counts for a sgRNA in the i^{th} sample. We assume X_i follows a binomial
 275 distribution, $X_i|p_i \sim Binomial(n_i, p_i)$, where n_i is the total number of mapped reads in sample i . To
 276 combine the estimated \hat{p}_i across multiple samples of the same treatment group, we proposed a linear
 277 model $p^A = \sum w_i p_i$, where i is the index for samples and w is the weight vector for samples in group A.
 278 Baggerly et al. (2003) proved that as long as $w^T \mathbf{1} = 1$, the expectation of $E(p^A)$ is unbiased. The value of
 279 w is estimated through gradient descent methods by minimizing the variance on p^A . Baggerly et al.

280 (2003) showed that $w_i \propto \left[\frac{1}{\alpha+\beta} + \frac{1}{n_i} \right]^{-1}$.

281 CB² performs the sgRNA-level differential analysis between two groups using a Student t -test like
 282 statistic (Baggerly et al. 2003):

$$t = \frac{p_B - p_A}{\sqrt{V_B + V_A}}$$

283 where p_A and p_B are the proportions of sgRNA, and V_A and V_B are the group variances of sgRNA, for
 284 groups A and B, respectively. Test statistic t represents the strength of the difference of sgRNA
 285 abundance between groups A and B. In other words, a large positive t -value indicates that the quantity of
 286 sgRNA in group B is greater than in group A, and a large negative t -value indicates that the quantity of
 287 sgRNA in group B is less than that in group A.

288 The variance is estimated by

$$\hat{v} = \max \left[\frac{\sum w_i^2 \hat{p}_i^2 - (\sum w_i^2) \hat{p}^2}{1 - (\sum w_i^2)}, \frac{\sum X_i \left(1 - \frac{\sum X_i}{\sum n_i}\right)}{\sum n_i} \right].$$

289 To measure the statistical significance of the difference, we approximate the p -value of a given t in a
 290 Student's t -distribution with a degree of freedom (df) defined by

$$df = \frac{(V_A + V_B)^2}{\frac{V_A^2}{n_A - 1} + \frac{V_B^2}{n_B - 1}},$$

291 where n_A and n_B are the numbers of replicates in groups A and B.

292 **sgRNA p -value aggregation for gene-level statistics**

293 Because multiple significant sgRNAs targeting the same gene hold greater biological significance than a
 294 single significant sgRNA, we must aggregate p -values to increase confidence in target identification. To
 295 do so, we combine p -values of sgRNAs for a target gene using Fisher's method (Fisher 1925) to assess
 296 overall differences at the gene level. The combined chi-square statistical test is used:

$$\chi_{2k}^2 \sim -2 \sum_{j=1}^k \ln(p_j),$$

297 where k is the number of sgRNAs targeting a gene in the screen and p_j is the p -value of j -th sgRNA for
 298 the gene. χ^2 follows a chi-squared distribution with $2k$ degrees of freedom. To correct for multiple
 299 hypothesis testing, we adapted the Benjamini-Hochberg procedure (Benjamini and Hochberg 1995) to
 300 estimate the False Discovery Rate (FDR).

301 **Gene-level statistics benchmarking on existing methods**

302 We used three different CRISPRn/CRISPRi pooled screen datasets (Evers et al. 2016) (RT112 and
303 UMUC3 cell line screens with CRISPR; RT112 cell line screen with CRISPRi) and two different
304 genome-wide CRISPRn/CRISPRi pooled screen datasets (Sanson et al. 2018) (A375 cell line screens),
305 which provide ground-truth labels of essentiality for each gene. With those screening datasets and labels,
306 we benchmarked the accuracy of essential gene detection by CB² with eight other published methods
307 (HiTSelect, MAGeCK, PBNPA, PinAPL-Py, RIGER, RSA, ScreenBEAM, and sgrSEA). We computed
308 the False Discovery Rate (FDR) for each gene from each method in the benchmark and set five different
309 levels of FDR cut-off (0.1, 0.05, 0.01, 0.005, 0.001) for essential gene classification. For example, if we
310 set FDR cut-off to 0.1, then a gene is predicted to be essential in the cell line if the FDR of the gene falls
311 below the cut-off value. We calculated recall (a recall value close to 1 indicates a prediction with a low
312 false negative rate), precision (a precision value close to 1 indicates a prediction with a low false positive
313 rate), and F-measure (the harmonic mean of precision and recall) of all the methods at each FDR level. To
314 allow each method to archive its best performance, we tuned parameters for the five parameter-tunable
315 methods – MAGeCK (permutation parameter for RRA test), PBNPA (no.sim parameter), RIGER (alpha
316 parameter), sgrSEA (multiplier parameter) and ScreenBEAM (burnin parameter) on CRISPRn-RT112
317 dataset. The F-measure was used as a measure for the parameter tuning. Most of the methods showed
318 robust performance regardless of the varied parameters, except sgrSEA and RIGER (Supplemental
319 Figure 17). Therefore, we used the default parameter for MAGeCK, PBNPA, and ScreenBEAM for other
320 datasets and used an optimized parameter for sgrSEA and RIGER. We also calculated the AUC (Area
321 Under the Curve) of the Precision and Recall curves and Receiver Operating Characteristic curves of all
322 the methods with FDR values. All of the data and scripts for the benchmarking are available at
323 <https://github.com/hyunhwaj/CB2-Experiments>. In addition, the data and the scripts can be found in the
324 Supplemental Material. Parameters used in these experiments are described below.

325 **CB²**

326 CB² R package (<https://CRAN.R-project.org/package=CB2>) were used in the benchmarking (R Core
327 Team 2019). Benchmarking of CB² was performed without parameter tuning since CB² is parameter-free.
328 FDR values for negative changes between two different groups from CB² statistical analysis were used
329 for benchmarking.

330 **HiTSelect**

331 We ran HiTSelect MATLAB package (<https://github.com/diazlab/HiTSelect>). Normalization by
332 Sequencing depth option was selected for benchmarking.

333 MAGeCK

334 MAGeCK version 0.5.8 was used for benchmarking. We ran MAGeCK with the ‘mageck test’ command
335 with the following parameters: ‘--norm-method’ ‘median’ and ‘--adjust-method’ ‘fdr’. We performed 100
336 permutations for the modified robust ranking aggregation (RRA) algorithm to estimate the gene-level
337 statistics on the benchmark datasets.

338 ScreenBEAM

339 ScreenBEAM R package (version 1.0.0, <https://github.com/jyyu/ScreenBEAM>) was used for
340 benchmarking ‘data.type’ parameter was set as ‘NGS,’ and ‘do.normalization’ was set as TRUE, ‘nitt’
341 and ‘burnin’ parameters for Bayesian computing were set at 15000 and 5000. ScreenBEAM does not
342 provide the one-sided p -value for negative selection, so for the FDR comparison with other methods, we
343 changed the FDR of a gene to 1 if the β of the gene is greater than 0.

344 PinAPL-Py

345 We used the PinAPL-Py website (<http://pinapl-py.ucsd.edu>) to perform the benchmarking. For the
346 sgRNA read counting, we used ‘GGCTTTATATATCTTGTGGAAAGGACGAAACACCG,
347 GCTTTATATATCTTGTGGAAAGGACGAAACACCG,’ and
348 ‘CTTTATATATCTTGTGGAAAGGACGAAACACCG,’ were used for ‘seq_5_end’ parameters of
349 ‘CRISPRn-RT112’, ‘CRISPRn-UMUC3’, and ‘CRISPRi-RT112’ datasets. We used CPM normalization
350 and set the GeneMetric parameter as ‘aRRA’ to perform a modified robust ranking aggregation (RRA).
351 We used the combined FDR values for each gene in the benchmarking.

352 RIGER

353 We used the Java implementation of RIGER (version 2.0, <https://github.com/broadinstitute/rigerj>) to
354 perform the benchmarking. We set the ‘alpha’ parameter at 0.1 on (Evers et al. 2016) datasets and at 1.0
355 on (Sanson et al. 2018) datasets. \log_2 fold-change values calculated by CB^2 were used as an input of
356 RIGER.

357 RSA

358 We used the Python implementation of RSA (version 1.9, <https://admin-ext.gnf.org/publications/RSA/>).
359 \log_2 fold-change values calculated by CB^2 were used as an input of RSA.

360 sgRSEA

361 sgRSEA R package (version 0.1, <https://cran.r-project.org/web/packages/sgRSEA/>) was used for
362 benchmarking. we set the multiplier at 30.

363 **PBNPA**

364 PBNPA R PACKAGE (version 0.0.2, <https://cran.r-project.org/web/packages/PBNPA/>) was used for
365 benchmarking. We set the sim.no parameter at 10.

366 **Specificity measure at sgRNA level**

367 The specificity of detecting true negative from the negative control sgRNAs is measured using
368 $\frac{\sum_{i=1}^N 1_A(p_i < \theta)}{N}$, where N is the number of non-targeting sgRNA, p_i is the estimated p-value of the i-th
369 sgRNA, and θ is the p-value threshold, and 1_A is the indicator function.

370 **Algorithm for quantifying sgRNA abundance**

371 **Previous sgRNA abundance quantification methods**

372 Recently published tools for CRISPR pooled screen analysis, including CRISPRcloud (Jeong et al. 2017),
373 MAGeCK (Li et al. 2014), CRISPRAnalyzeR (Winter et al. 2017), and PinAPL-Py (Spahn et al. 2017)
374 provide different methods for estimating the abundance of sgRNAs in each sample from pooled libraries.
375 In most cases, input data consist of raw FASTQ-format sequencing result files.

376 CRISPRcloud was the first tool to offer an online user-defined, light-weight quantification method which
377 proceeds on the user-client side. In contrast, CRISPRAnalyzeR and PinAPL-Py run their quantification
378 methods on the server-side. As a result, CRISPRcloud minimized information passed through the Internet
379 by transferring only the processed count matrix to the cloud storage.

380 However, as pointed out by Spahn et al. (2017), CRISPRcloud quantification algorithm can produce
381 erroneous mapping results if the sgRNA sequences are staggered, because CRISPRcloud extracts sgRNA
382 sequence for each read at a fixed location. Another limitation of CRISPRcloud is the fact that the user
383 must decide where the extraction site is. Nevertheless, CRISPRcloud did not require tuning and was thus
384 arguably more user-friendly than other tools. For instance, in PinAPL-Py, users need to set many tuning
385 parameters for sgRNA quantification: adapter error rate parameter for trimming, matching and ambiguity
386 thresholds and parameters for alignment seed for the Bowtie alignment (Spahn et al. 2017).

387 We engineered CB² to address precisely these issues. As a result, users no longer need to perform
388 complicated parameter tuning for the sgRNA abundance quantification; one must simply provide the
389 input files to CB².

390 **The binary representation of sgRNA sequence lowers the cost of computation**

391 We used a binary representation for sgRNA sequence. This approach is memory-efficient and improves
392 the user experience at the client-side (Melsted and Pritchard 2011). It only needs $\max(K, 2M)$ bits to
393 store an sgRNA-sequence, where M is the length of the sequence and K is the bit size to store a primitive
394 integer in the machine (usually 64 bits) because we only need two bits to save a nucleotide (i.e., ‘A’ is
395 ‘00’, ‘C’ is ‘01’, ‘G’ is ‘11’, and ‘T’ is ‘10’). The memory size is about half of that required for storing a
396 character string of the sequence, i.e., 160 bits are needed to store a 20nt sgRNA sequence. Another benefit
397 of binary representation is that it lowers the time complexity for the shift operator when comparing all k -
398 mers of an sgRNA read using a sliding window. This is an essential function for the quantification
399 algorithm in CB². Compared to the string shift operator functions, such as string copy, substring
400 extraction, and concatenation, the binary representation produces significantly shorter running times.
401 Algorithm 2 in Supplemental Methods illustrates how a sgRNA library converted to a bit sequence and
402 store the converted sequence into a hash table.

403 **Sliding window-based algorithm gives a high-resolution quantification with comparable** 404 **running time**

405 With a binary representation, we run the quantification algorithm as follows: First, we build a hash table
406 for the reference library, with each key of the library in the hash table converted to the binary
407 representation. Second, for each read, we scan the sequence of the read from 5' to 3' with the sliding
408 window. In the i -th iteration, the sliding window contains a substring of the read sequence from i to
409 $i + k - 1$, where k is the length of the sgRNAs. The substring is also converted to a binary sequence, and
410 the hash table is quickly checked to see if the sequence in the sliding window exists in the reference
411 library. If the sequence is found in the hash table, then the count of the sequence is increased by one, and
412 then the algorithm proceeds to the next read. Otherwise, it moves to $i + 1$ -th iteration and the bit-shift
413 method will be applied to take the next sliding window. Algorithm 1 in Supplemental Methods represents
414 a procedure of the Sliding window-based algorithm.

415 For the case of a reverse complement sequenced sample, the entire procedure is repeated on the reverse
416 complement reference sgRNA library and scanning the read from 3' to 5'. After both assays are
417 performed (5' to 3' and 3' to 5' with the reverse complement reference sequences), mapping results
418 between both sequences are compared. The one with a larger count corresponds to the correct sequence
419 mapping. We compared the mappability of CB² to those of MAGeCK (Li et al. 2015) and PinAPL-Py
420 (Spahn et al. 2017) across multiple datasets from previous studies (Figure 5).

421 **Acknowledgments**

422 This work has been supported by National Institute of General Medical Sciences R01-GM120033,
 423 National Science Foundation - Division of Mathematical Sciences DMS-1263932, Cancer Prevention
 424 Research Institute of Texas RP170387, Houston Endowment, and Chao Family Foundation (Z.L.),
 425 Huffington Foundation, Howard Hughes Medical Institute (H.Y.Z.), and the Parkinson's Foundation
 426 Stanley Fahn Junior Faculty Award PF-JFA-1762 (M.W.C.R.). We thank V. L. Brandt for editing the
 427 manuscript.

428 **Author contributions**

429 H-H.J., M.W.C.R., H.Y.Z., and Z.L. designed the study. H-H.J. and S.Y.K. implemented the CB²
 430 software. H-H.J. and Z.L. performed analysis. Z.L. supervised the project. H-H.J., M.W.C.R., and Z.L.
 431 wrote the manuscript with input from all the authors.

432 **Competing interests**

433 The authors declare no competing interests.

434 **References**

- 435 Aguirre AJ, Meyers RM, Weir BA, Vazquez F, Zhang C-Z, Ben-David U, Cook A, Ha G, Harrington
 436 WF, Doshi MB, et al. 2016. Genomic Copy Number Dictates a Gene-Independent Cell Response
 437 to CRISPR/Cas9 Targeting. *Cancer Discovery*.
- 438 Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome biology* **11**:
 439 R106.
- 440 Baggerly KA, Deng L, Morris JS, Aldaz CM. 2003. Differential expression in SAGE: accounting for
 441 normal between-library variation. *Bioinformatics* **19**: 1477–1483.
- 442 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful approach
 443 to multiple testing. *Journal of the Royal Statistical Society* **57**: 289–300.
- 444 DeJesus R, Moretti F, McAllister G, Wang Z, Bergman P, Liu S, Frias E, Alford J, Reece-Hoyes JS,
 445 Lindeman A, et al. 2016. Functional CRISPR screening identifies the ufmylation pathway as a
 446 regulator of SQSTM1/p62. *eLife* **5**.
- 447 Diaz AA, Qin H, Ramalho-Santos M, Song JS. 2015. HiTSelect: a comprehensive tool for high-
 448 complexity-pooled screen analysis. *Nucleic acids research* **43**: e16.
- 449 Doench JG. 2017. Am I ready for CRISPR? A user's guide to genetic screens. *Nature Reviews Genetics*
 450 **19**: 67–80.

- 451 Dutta B, Azhir A, Merino L-H, Guo Y, Revanur S, Madhamshettiwar PB, Germain RN, Smith JA,
452 Simpson KJ, Martin SE, et al. 2016. An interactive web-based application for Comprehensive
453 Analysis of RNAi-screen Data. *Nature communications* **7**: 10578.
- 454 Evers B, Jastrzebski K, Heijmans JPM, Grenrum W, Beijersbergen RL, Bernards R. 2016. CRISPR
455 knockout screening outperforms shRNA and CRISPRi in identifying essential genes. *Nature*
456 *biotechnology* **34**: 11–14.
- 457 Fisher RA. 1925. *Statistical methods for research workers*. Genesis Publishing Pvt Ltd.
- 458 Gilbert LA, Horlbeck MA, Adamson B, Villalta JE, Chen Y, Whitehead EH, Guimaraes C, Panning B,
459 Ploegh HL, Bassik MC, et al. 2014. Genome-Scale CRISPR-Mediated Control of Gene
460 Repression and Activation. *Cell* **159**: 647–61.
- 461 Golden RJ, Chen B, Li T, Braun J, Manjunath H, Chen X, Wu J, Schmid V, Chang T-C, Kopp F, et al.
462 2017. An Argonaute phosphorylation cycle promotes microRNA-mediated silencing. *Nature* **542**:
463 197–202.
- 464 Hart T, Brown KR, Sircoulomb F, Rottapel R, Moffat J. 2014. Measuring error rates in genomic
465 perturbation screens: gold standards for human functional genomics. *Molecular systems biology*
466 **10**: 733.
- 467 Hart T, Chandrashekar M, Aregger M, Steinhart Z, Brown KR, MacLeod G, Mis M, Zimmermann M,
468 Fradet-Turcotte A, Sun S, et al. 2015. High-Resolution CRISPR Screens Reveal Fitness Genes
469 and Genotype-Specific Cancer Liabilities. *Cell* **163**: 1515–1526.
- 470 Horlbeck MA, Gilbert LA, Villalta JE, Adamson B, Pak RA, Chen Y, Fields AP, Park CY, Corn JE,
471 Kampmann M, et al. 2016. Compact and highly active next-generation libraries for CRISPR-
472 mediated gene repression and activation. *Elife* **5**: e19760.
- 473 Hu G, Luo J. 2012. A primer on using pooled shRNA libraries for functional genomic screens. *Acta*
474 *Biochimica et Biophysica Sinica* **44**: 103–112.
- 475 Iorio F, Behan FM, Gonçalves E, Bhosle SG, Chen E, Shepherd R, Beaver C, Ansari R, Pooley R,
476 Wilkinson P, et al. 2018. Unsupervised correction of gene-independent cell responses to
477 CRISPR-Cas9 targeting. *BMC Genomics* **19**: 604.
- 478 Jeong HH, Kim SY, Rousseaux MW, Zoghbi HY, Liu Z. 2017. CRISPRcloud: A secure cloud-based
479 pipeline for CRISPR pooled screen deconvolution. *Bioinformatics* **33**: 2963–2965.
- 480 Jia G, Wang X, Xiao G. 2017. A permutation-based non-parametric analysis of CRISPR screen data.
481 *BMC genomics* **18**: 545.
- 482 Koike-Yusa H, Li Y, Tan E-P, Velasco-Herrera MDC, Yusa K. 2014. Genome-wide recessive genetic
483 screening in mammalian cells with a lentiviral CRISPR-guide RNA library. *Nature biotechnology*
484 **32**: 267–73.
- 485 König R, Chiang C, Tu BP, Yan SF, DeJesus PD, Romero A, Bergauer T, Orth A, Krueger U, Zhou Y, et
486 al. 2007. A probability-based approach for the analysis of large-scale RNAi screens. *Nature*
487 *Methods* **4**: 847–849.

- 488 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature methods* **9**: 357.
- 489 Li M, Yu JSL, Tilgner K, Ong SH, Koike-Yusa H, Yusa K. 2018. Genome-wide CRISPR-KO Screen
490 Uncovers mTORC1-Mediated Gsk3 Regulation in Naive Pluripotency Maintenance and
491 Dissolution. *Cell Reports*.
- 492 Li W, Köster J, Xu H, Chen C-H, Xiao T, Liu JS, Brown M, Liu XS. 2015. Quality control, modeling,
493 and visualization of CRISPR screens with MAGeCK-VISPR. *Genome biology* **16**: 281.
- 494 Li W, Xu H, Xiao T, Cong L, Love MI, Zhang F, Irizarry RA, Liu JS, Brown M, Liu XS. 2014.
495 MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9
496 knockout screens. *Genome biology* **15**: 554.
- 497 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-seq
498 data with DESeq2. *Genome biology* **15**: 550.
- 499 Luo B, Cheung HW, Subramanian A, Sharifnia T, Okamoto M, Yang X, Hinkle G, Boehm JS,
500 Beroukhim R, Weir BA, et al. 2008. Highly parallel identification of essential genes in cancer
501 cells. *Proceedings of the National Academy of Sciences* **105**: 20380–20385.
- 502 Luo J, Emanuele MJ, Li D, Creighton CJ, Schlabach MR, Westbrook TF, Wong K-K, Elledge SJ. 2009.
503 A Genome-wide RNAi Screen Identifies Multiple Synthetic Lethal Interactions with the Ras
504 Oncogene. *Cell* **137**: 835–848.
- 505 Marioni JC, Mason CE, Mane SM, Stephens M, Gilad Y. 2008. RNA-seq: an assessment of technical
506 reproducibility and comparison with gene expression arrays. *Genome research* **18**: 1509–1517.
- 507 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
508 *EMBnet.journal* **17**: 10.
- 509 Melsted P, Pritchard JK. 2011. Efficient counting of k-mers in DNA sequences using a bloom filter. *BMC*
510 *Bioinformatics* **12**: 333.
- 511 Mohr SE, Smith JA, Shamu CE, Neumüller RA, Perrimon N. 2014. RNAi screening comes of age:
512 improved techniques and complementary approaches. *Nature reviews Molecular cell biology* **15**:
513 591–600.
- 514 Noh J. 2015. *sgRSEA: Enrichment Analysis of CRISPR/Cas9 Knockout Screen Data*. [https://cran.r-](https://cran.r-project.org/web/packages/sgRSEA/)
515 [project.org/web/packages/sgRSEA/](https://cran.r-project.org/web/packages/sgRSEA/).
- 516 Paddison PJ, Silva JM, Conklin DS, Schlabach M, Li M, Aruleba S, Baliya V, O’Shaughnessy A, Gnoj L,
517 Scobie K, et al. 2004. A resource for large-scale RNA-interference-based screens in mammals.
518 *Nature* **428**: 427–431.
- 519 Park J, Al-Ramahi I, Tan Q, Mollema N, Diaz-Garcia JR, Gallego-Flores T, Lu H-C, Lagalwar S, Duvick
520 L, Kang H, et al. 2013. RAS–MAPK–MSK1 pathway modulates ataxin 1 protein levels and
521 toxicity in SCA1. *Nature* **498**: 325–331.
- 522 Parnas O, Jovanovic M, Eisenhaure TM, Herbst RH, Dixit A, Ye CJ, Przybylski D, Platt RJ, Tirosh I,
523 Sanjana NE, et al. 2015. A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect
524 Regulatory Networks. *Cell* **162**: 675–86.

- 525 R Core Team. 2019. *R: A Language and Environment for Statistical Computing*. R Foundation for
526 Statistical Computing, Vienna, Austria <https://www.R-project.org>.
- 527 Sanjana NE, Shalem O, Zhang F. 2014. Improved vectors and genome-wide libraries for CRISPR
528 screening. *Nature methods* **11**: 783.
- 529 Sanson KR, Hanna RE, Hegde M, Donovan KF, Strand C, Sullender ME, Vaimberg EW, Goodale A,
530 Root DE, Piccioni F, et al. 2018. Optimized libraries for CRISPR-Cas9 genetic screens with
531 multiple modalities. *Nature communications* **9**: 5416.
- 532 Schlabach MR, Luo J, Solimini NL, Hu G, Xu Q, Li MZ, Zhao Z, Smogorzewska A, Sowa ME, Ang XL,
533 et al. 2008. Cancer proliferation gene discovery through functional genomics. *Science (New York,
534 NY)* **319**: 620–4.
- 535 Shalem O, Sanjana NE, Hartenian E, Shi X, Scott DA, Mikkelsen TS, Heckl D, Ebert BL, Root DE,
536 Doench JG, et al. 2014. Genome-Scale CRISPR-Cas9 Knockout Screening in Human Cells.
537 *Science* **343**: 84–87.
- 538 Shao DD, Tsherniak A, Gopal S, Weir BA, Tamayo P, Stransky N, Schumacher SE, Zack TI, Beroukhim
539 R, Garraway LA, et al. 2013. ATARiS: computational quantification of gene suppression
540 phenotypes from multisample RNAi screens. *Genome research* **23**: 665–678.
- 541 Silva JM, Marran K, Parker JS, Silva J, Golding M, Schlabach MR, Elledge SJ, Hannon GJ, Chang K.
542 2008. Profiling Essential Genes in Human Mammary Cells by Multiplex RNAi Screening.
543 *Science* **319**: 617–620.
- 544 Simon MM, Moresco EMY, Bull KR, Kumar S, Mallon A-M, Beutler B, Potter PK. 2015. Current
545 strategies for mutation detection in phenotype-driven screens utilising next generation
546 sequencing. *Mammalian genome* □: *official journal of the International Mammalian Genome
547 Society* **26**: 486–500.
- 548 Spahn PN, Bath T, Weiss RJ, Kim J, Esko JD, Lewis NE, Harismendy O. 2017. PinAPL-Py: A
549 comprehensive web-application for the analysis of CRISPR/Cas9 screens. *Scientific reports* **7**:
550 15854.
- 551 Wang T, Wei JJ, Sabatini DM, Lander ES. 2014. Genetic screens in human cells using the CRISPR-Cas9
552 system. *Science (New York, NY)* **343**: 80–4.
- 553 Winter J, Schwering M, Pelz O, Rauscher B, Zhan T, Heigwer F, Boutros M. 2017. CRISPRAnalyzer:
554 Interactive analysis, annotation and documentation of pooled CRISPR screens. *bioRxiv*.
- 555 Yu J, Silva J, Califano A. 2016. ScreenBEAM: a novel meta-analysis algorithm for functional genomics
556 screens via Bayesian hierarchical modeling. *Bioinformatics (Oxford, England)* **32**: 260–7.
- 557 Zhou Y, Zhu S, Cai C, Yuan P, Li C, Huang Y, Wei W. 2014. High-throughput screening of a
558 CRISPR/Cas9 library for functional genomics in human cells. *Nature* **509**: 487–491.
- 559

560 Figure legends

561

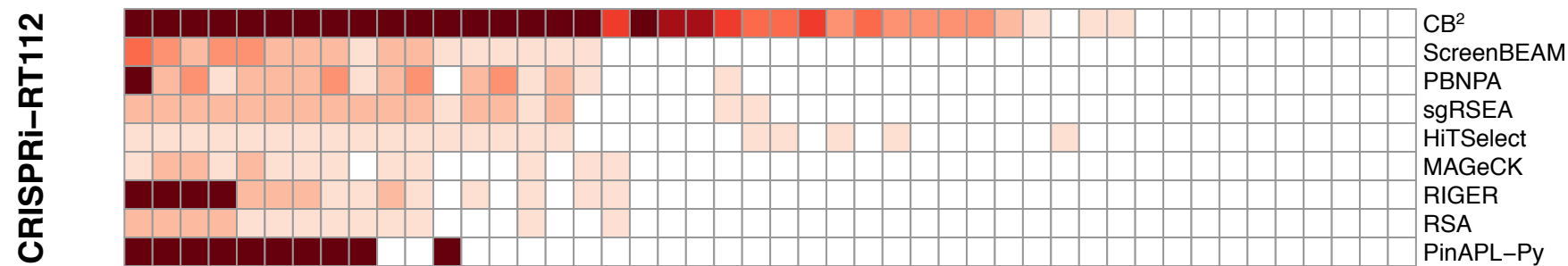
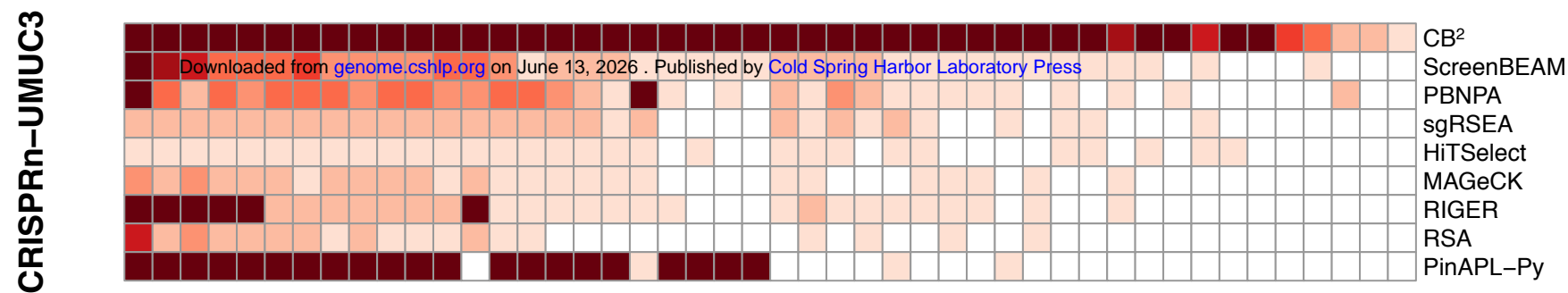
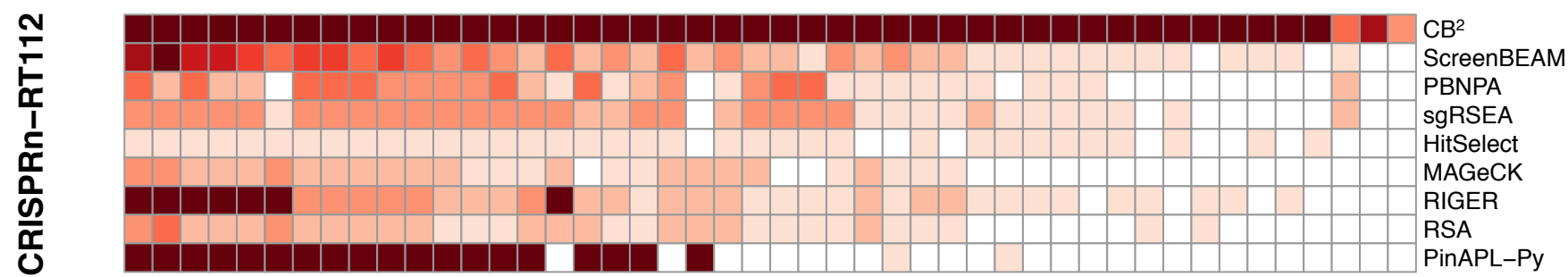
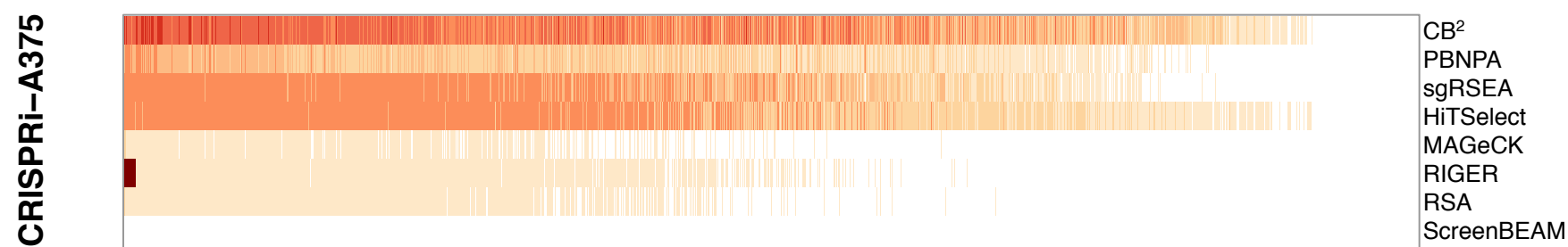
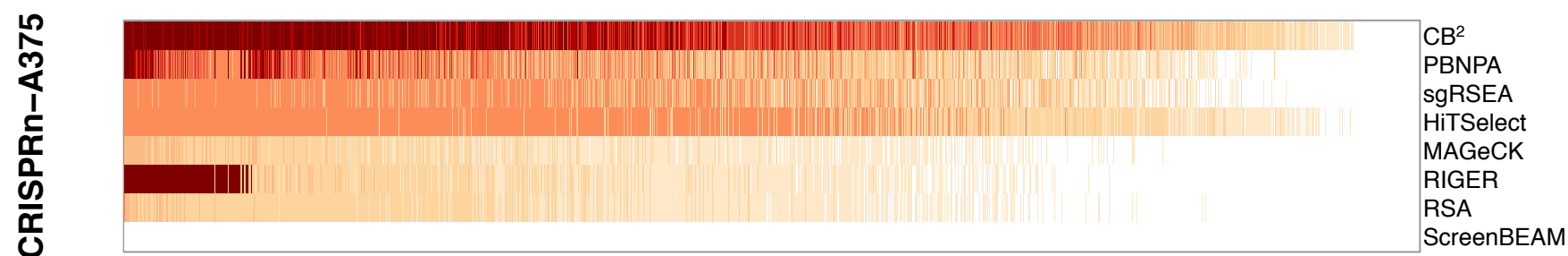
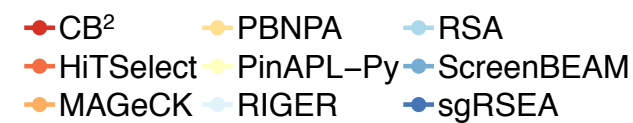
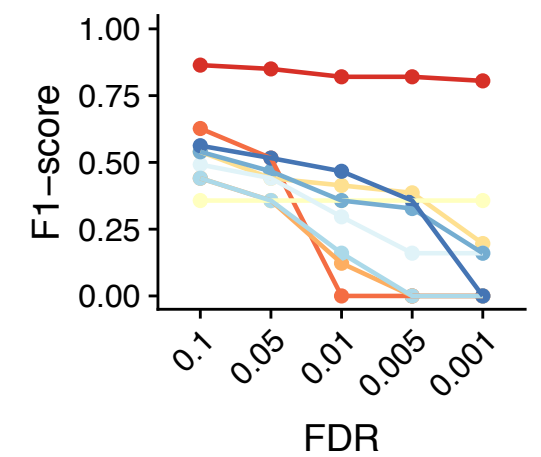
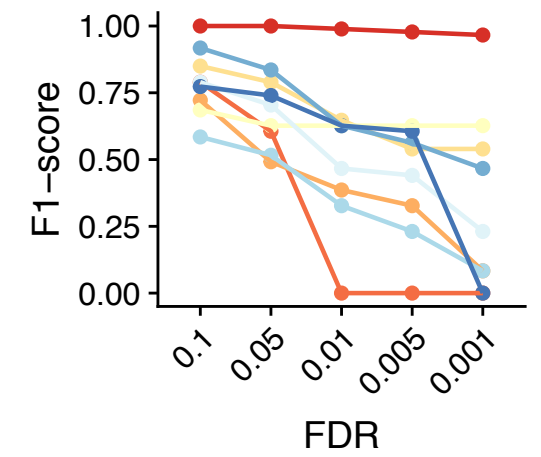
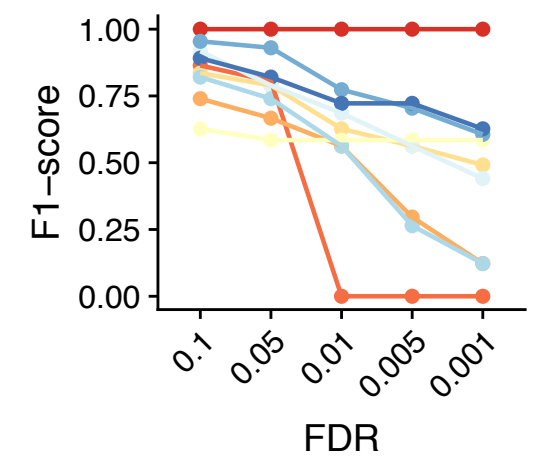
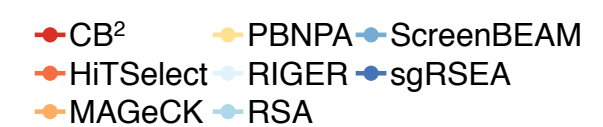
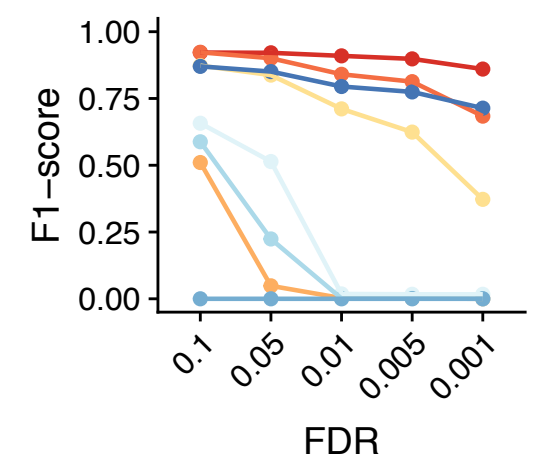
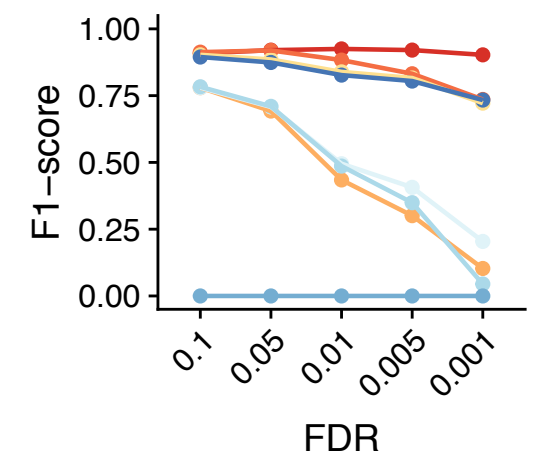
562 Figure 1. CB^2 offers robust target identification with high precision and recall. (A-B) Benchmark results using data
 563 from Evers et al. (2016). (A) Heatmaps illustrate FDRs of gene statistics from each of nine leading high-complexity
 564 pooled screen analysis tools. (B) F1-score measurements at different FDR cut-offs across all methods. At commonly
 565 used FDR cut-offs, CB^2 can identify most of the essential genes with high rates of precision and recall. (C-D) Same
 566 representation as in (A-B), using data from Sanson et al. (2018).

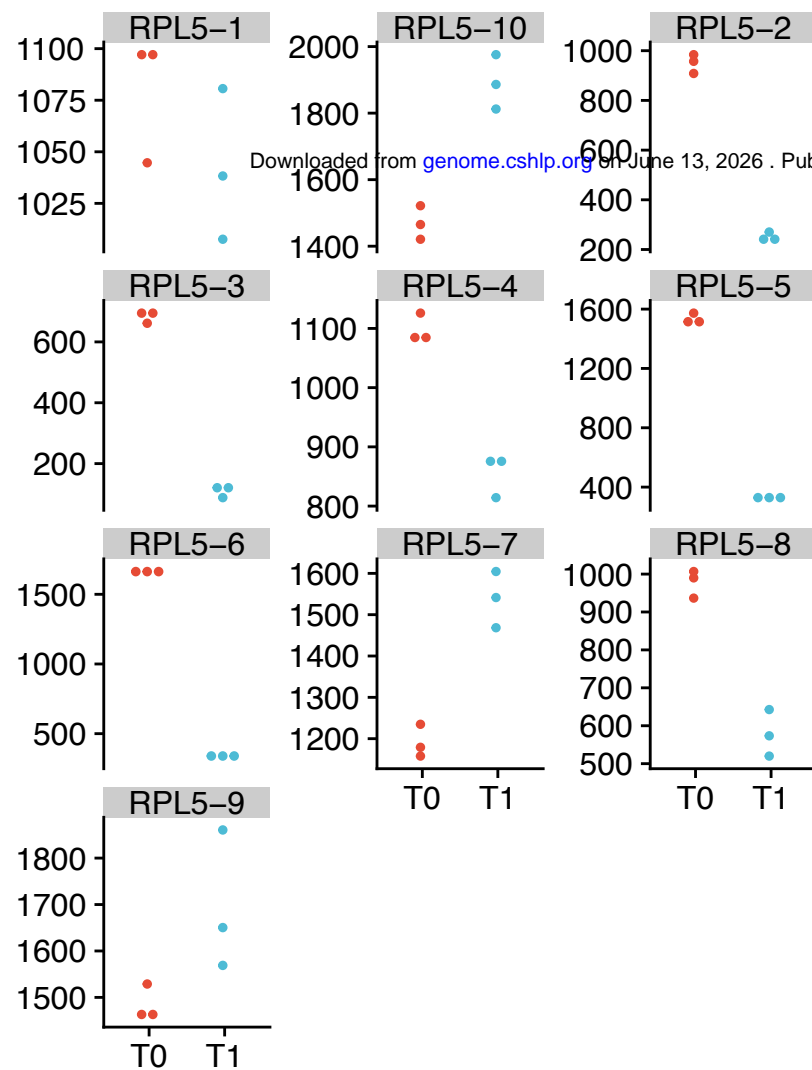
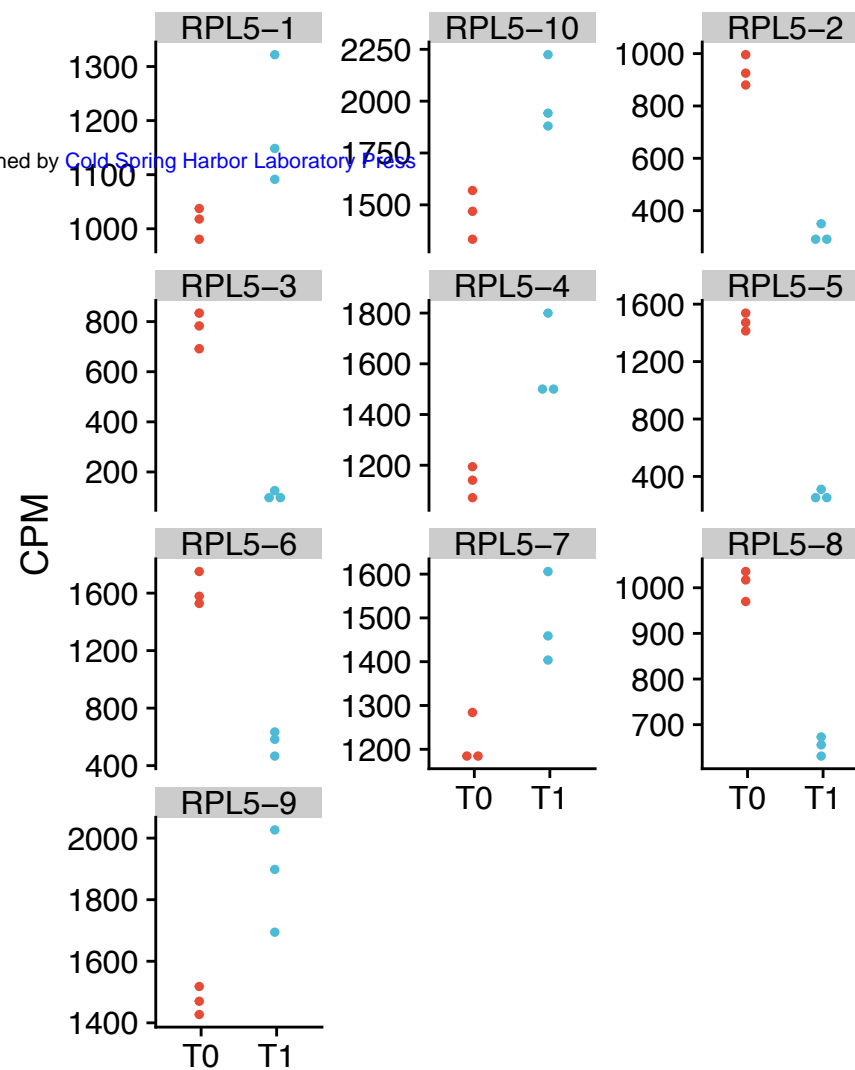
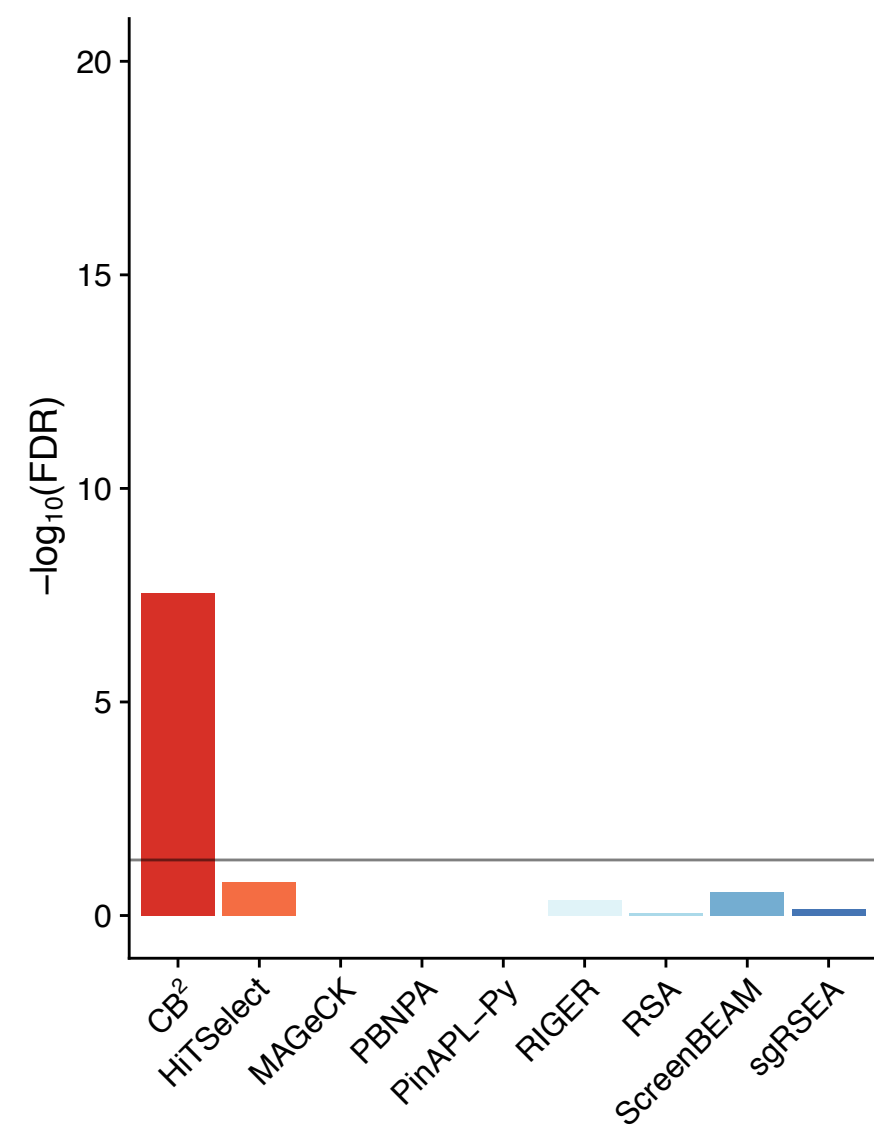
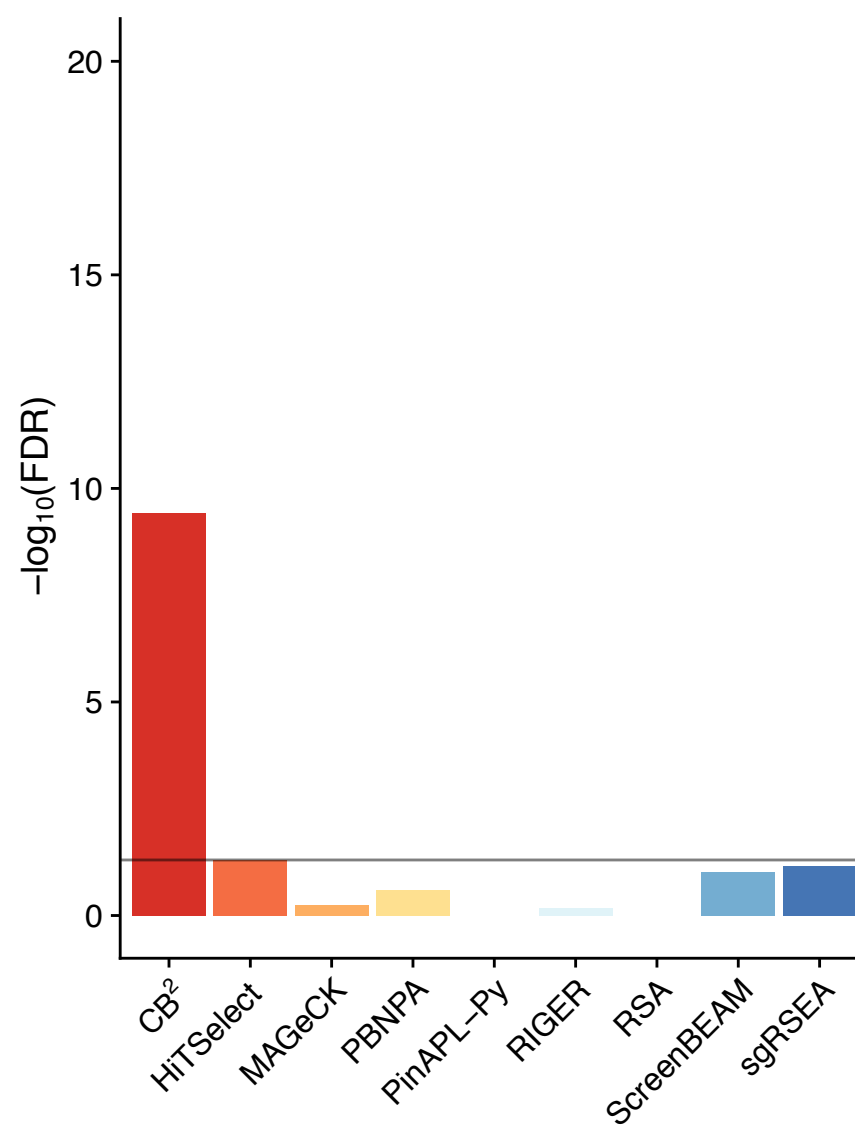
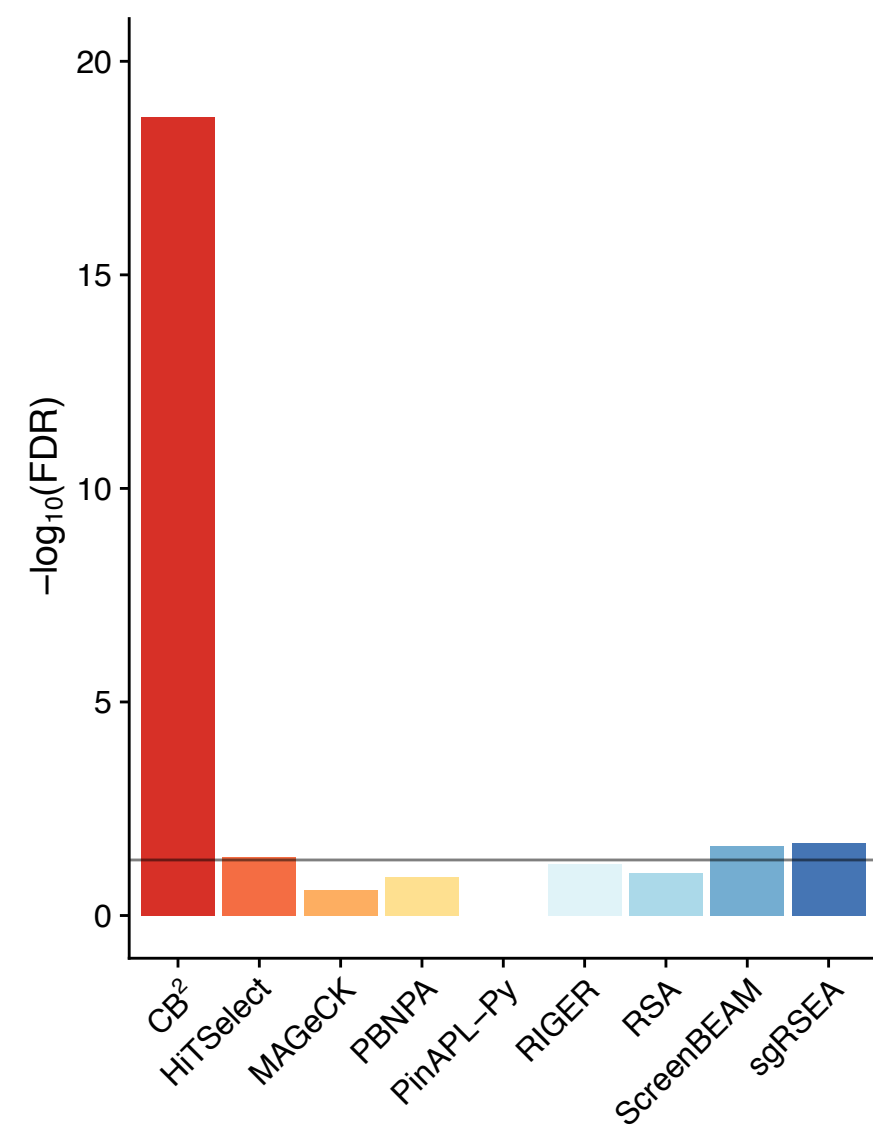
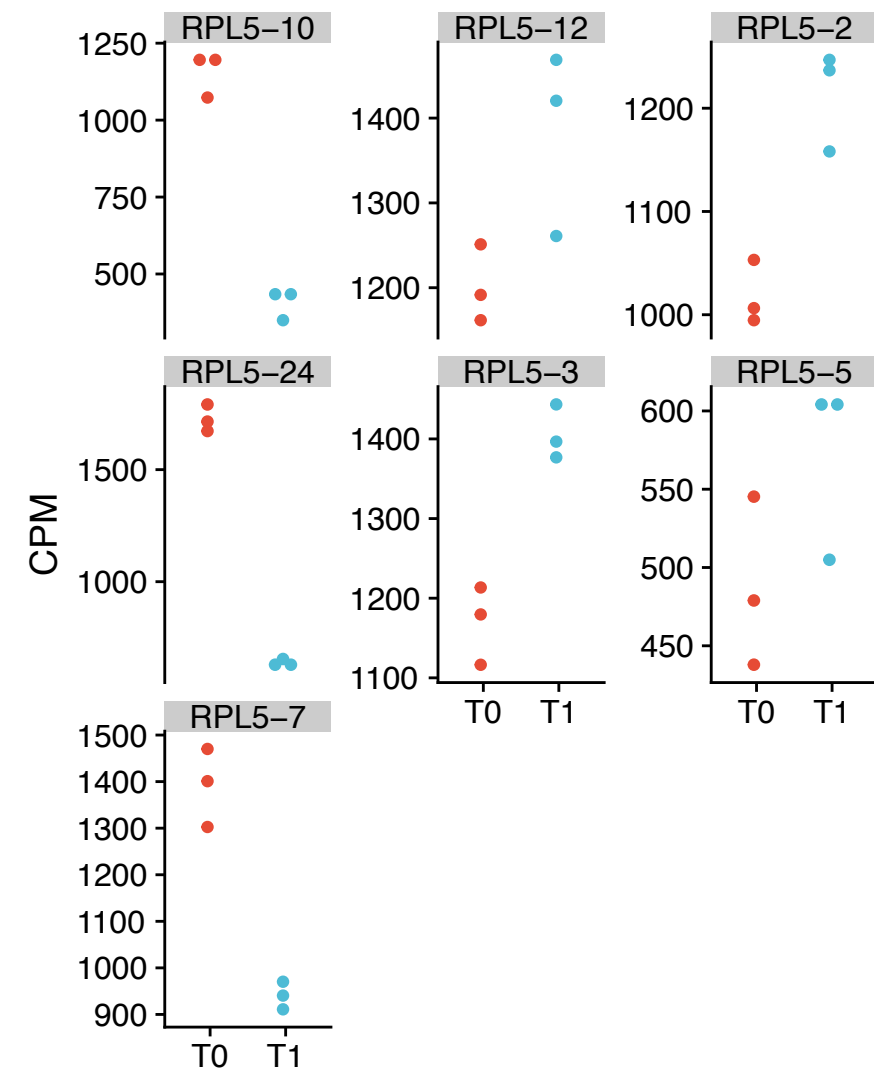
567 Figure 2. CB^2 detects essential genes missed by other leading methods: the case of *RPL5*. sgRNA quantification for
 568 *RPL5* in cell line (A) RT112, (B) UMUC3 using CRISPRn and (C) RT112 using the CRISPRi library. The top
 569 panels show CPM (count per million) of sgRNAs that target *RPL5* for each group (T_0 and T_1), and the bottom panels
 570 indicate the reported the FDR for *RPL5* in each screen across all the methods. A horizontal line at $FDR = 0.01$ is
 571 used as a threshold for statistical cutoff. CB^2 outperforms all other methods of identifying *RPL5* as an essential gene
 572 across all benchmark datasets.

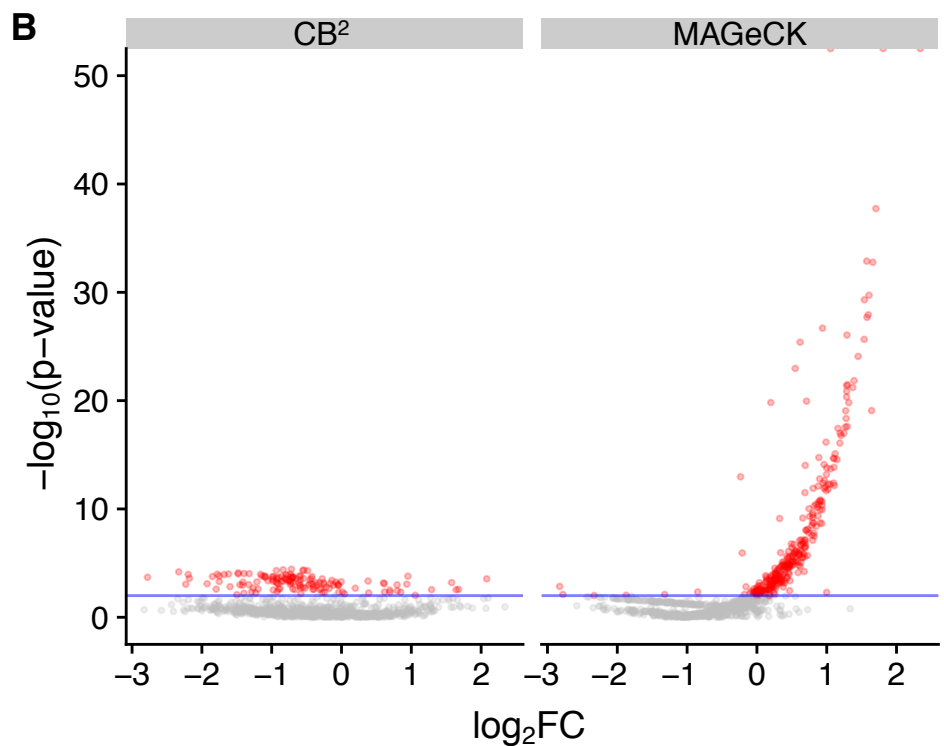
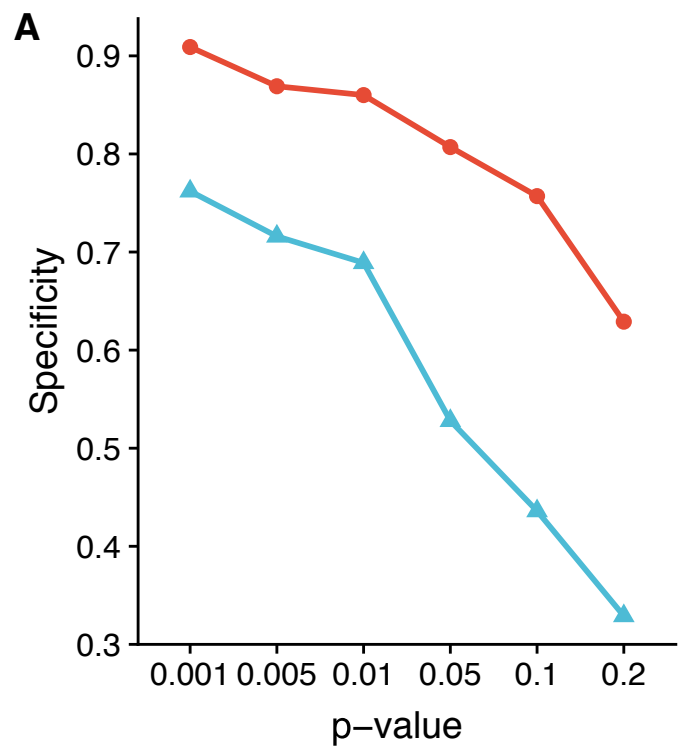
573 Figure 3. Comparison of False Positive Rate for non-targeting sgRNAs on (Parnas et al. 2015)'s screen data. (A)
 574 Specificity comparison between CB^2 and MAGeCK for the six different p-value thresholds. The y-axis indicates
 575 specificity, and the x-axis indicates the level of the p-value threshold for the specificity calculation. (B) Volcano
 576 plots of the p-value of non-targeting sgRNAs. The y-axis indicates the negative logarithm value of p-value, and the
 577 x-axis indicates the \log_2 value of fold-change. All of the data points are from negative control sgRNAs. False
 578 Positive were plotted in red. Horizontal blue lines at $p = 0.01$ indicates the threshold for statistical cutoff.

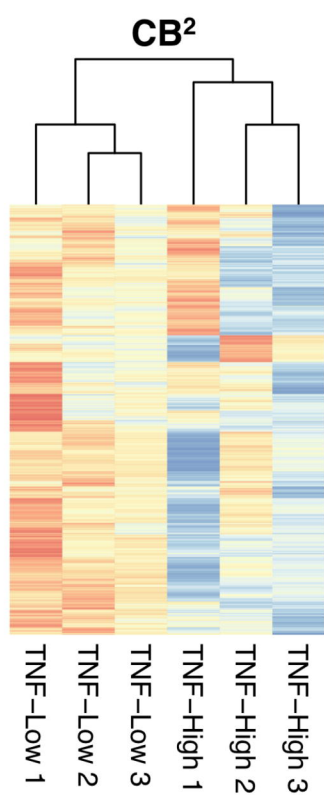
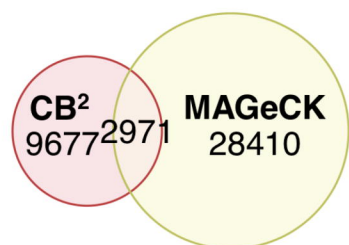
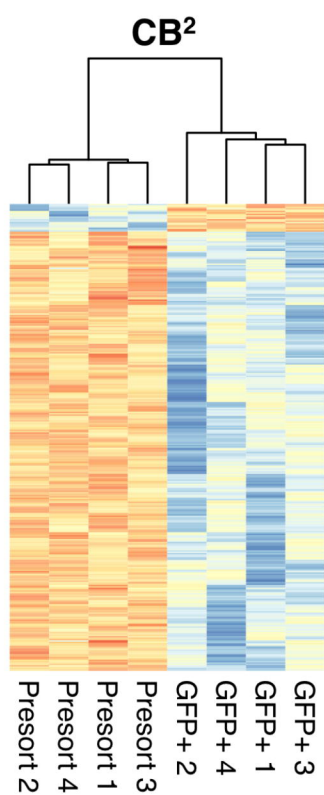
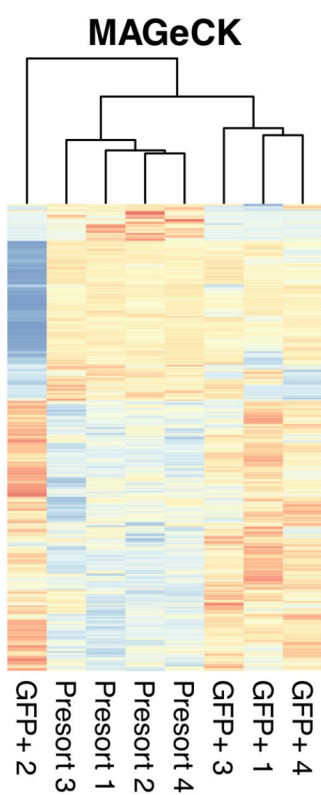
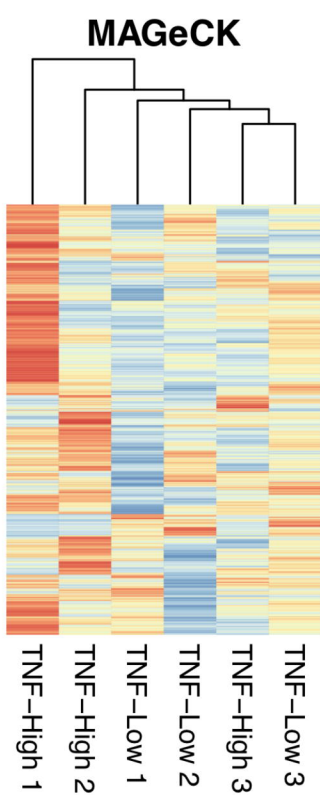
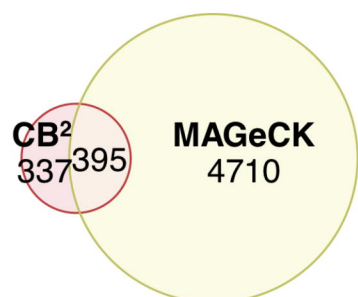
579 Figure 4. Discrepancy of sgRNA-level statistics between CB^2 and MAGeCK in two public CRISPR pooled screens.
 580 (A) Heatmaps of normalized read counts of the detected sgRNAs from the screen data in Parnas et al. (2015). Left:
 581 A heatmap of sgRNAs detected by CB^2 only. Right: A heatmap of sgRNAs detected by MAGeCK only. (B) Venn
 582 diagrams of sgRNAs detected by CB^2 and MAGeCK from the screen data in Parnas et al. (2015). (C-D) Same
 583 representations of (A-B) using data from Li et al. (2018).

584 Figure 5. CB^2 outperforms MAGeCK and PinAPL-Py in the percentage of mapped reads over six benchmark
 585 datasets. (A) Read mappability of CB^2 , MAGeCK, and PinAPL-Py across six different datasets. (B) Representative
 586 examples of reads that were not mapped by MAGeCK or PinAPL-Py. Adapters are highlighted with green, sgRNAs
 587 with a red. Yellow boxes show the predicted locations of sgRNAs by each method. Several parameters were used to
 588 optimize performances of PinAPL-Py.

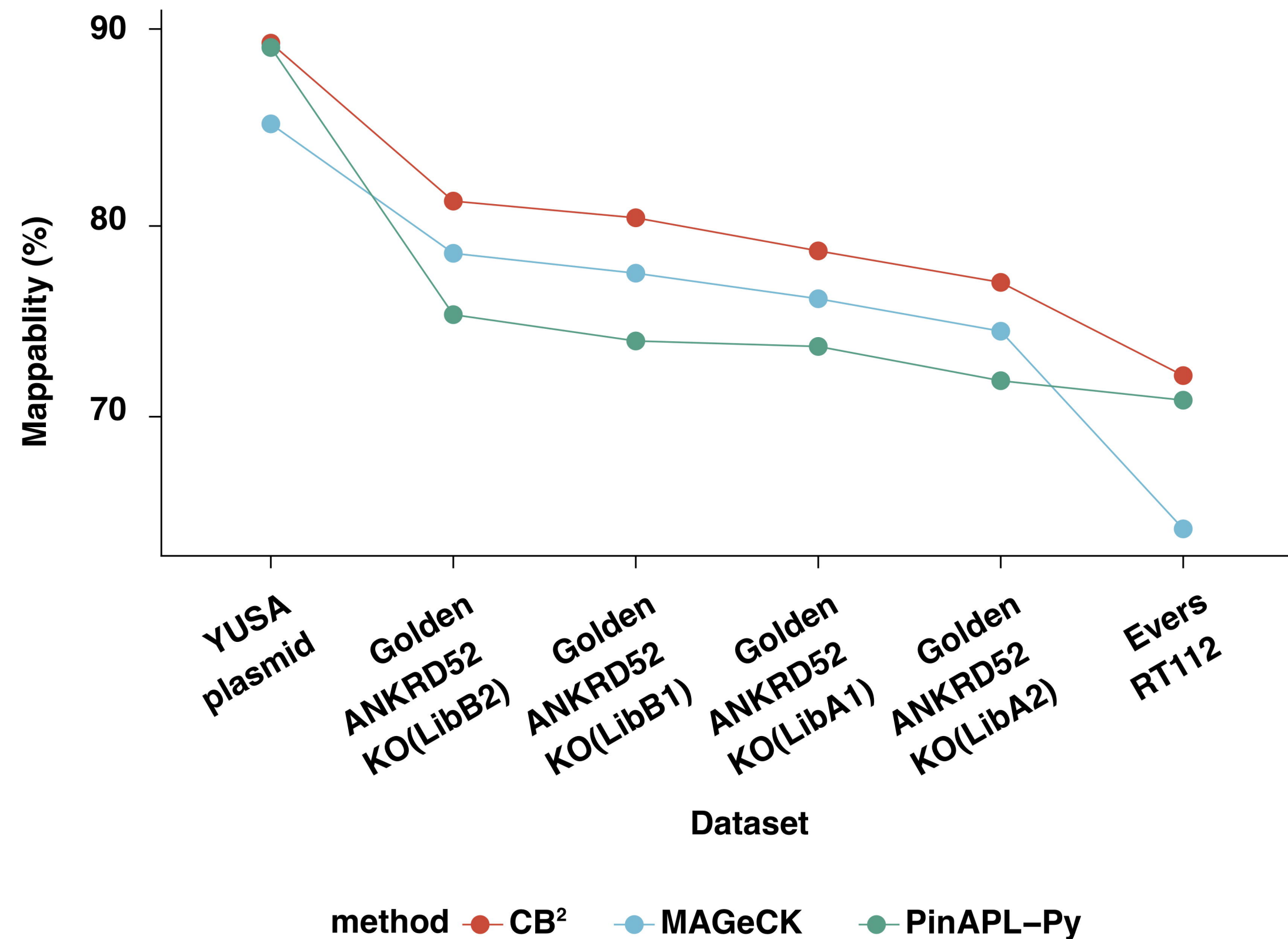
A $-\log_{10}(\text{FDR})$ **C** $-\log_{10}(\text{FDR})$ **B****D**

A**RPL5 (CRISPRn-RT112)****B****RPL5 (CRISPRn-UMUC3)****C****RPL5 (CRISPRi-RT112)**



A**B****C****D**

A



B

Extracted sgRNA sequence

Ground Truth

```
CGTGATGTCTTTATTGATCTTGTAGCCAGACCGCACCACCAGCTCGTGGACCTTGTGAGTTTA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAACGACCAGACACCACTGATTGCGTTTAAGA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAAACCACTCGGCAAACCAGGTCGTTTAAGAG
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAACGTTGGCCGGACTCTAGTGGCAGTTTAAGA
CGTGATGGCTTTATAATCTTGTGGAAAGGACGAACGAGACCCTGGTGAGCGTTGAGTTTAAGAG
```

CB²

```
CGTGATGTCTTTATTGATCTTGTAGCCAGACCGCACCACCAGCTCGTGGACCTTGTGAGTTTA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAACGACCAGACACCACTGATTGCGTTTAAGA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAAACCACTCGGCAAACCAGGTCGTTTAAGAG
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAACGTTGGCCGGACTCTAGTGGCAGTTTAAGA
CGTGATGGCTTTATAATCTTGTGGAAAGGACGAACGAGACCCTGGTGAGCGTTGAGTTTAAGAG
```

MAGeCK

```
CGTGATGTCTTTATTGATCTTGTAGCCAGACCGCACCACCAGCTCGTGGACCTTGTGAGTTTA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAACGACCAGACACCACTGATTGCGTTTAAGA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAAACCACTCGGCAAACCAGGTCGTTTAAGAG
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAACGTTGGCCGGACTCTAGTGGCAGTTTAAGA
CGTGATGGCTTTATAATCTTGTGGAAAGGACGAACGAGACCCTGGTGAGCGTTGAGTTTAAGAG
```

PinAPL-Py (e=0.1)

```
CGTGATGTCTTTATTGATCTTGTAGCCAGACCGCACCACCAGCTCGTGGACCTTGTGAGTTTA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAACGACCAGACACCACTGATTGCGTTTAAGA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAAACCACTCGGCAAACCAGGTCGTTTAAGAG
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAACGTTGGCCGGACTCTAGTGGCAGTTTAAGA
CGTGATGGCTTTATAATCTTGTGGAAAGGACGAACGAGACCCTGGTGAGCGTTGAGTTTAAGAG
```

PinAPL-Py (e=0.2)

```
CGTGATGTCTTTATTGATCTTGTAGCCAGACCGCACCACCAGCTCGTGGACCTTGTGAGTTTA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAACGACCAGACACCACTGATTGCGTTTAAGA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAAACCACTCGGCAAACCAGGTCGTTTAAGAG
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAACGTTGGCCGGACTCTAGTGGCAGTTTAAGA
CGTGATGGCTTTATAATCTTGTGGAAAGGACGAACGAGACCCTGGTGAGCGTTGAGTTTAAGAG
```

PinAPL-Py (e=0.3)

```
CGTGATGTCTTTATTGATCTTGTAGCCAGACCGCACCACCAGCTCGTGGACCTTGTGAGTTTA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAACGACCAGACACCACTGATTGCGTTTAAGA
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAAAACCACTCGGCAAACCAGGTCGTTTAAGAG
CGTGATGGCTTTATATATCTTGTGGAAAGGACGAACGTTGGCCGGACTCTAGTGGCAGTTTAAGA
CGTGATGGCTTTATAATCTTGTGGAAAGGACGAACGAGACCCTGGTGAGCGTTGAGTTTAAGAG
```

Table 1. Statistical models used by CB2 and existing methods. All of the methods were used in the target identification benchmarking.

Name	sgRNA-level statistics	gene-level statistics
CB ²	Beta-binomial distribution	Fisher's method
HitSelect (Diaz et al. 2015)	Poisson distribution (active number of sgRNAs)	Stochastic multi-objective ranking method for gene-level statistics
MAGeCK (Li et al. 2015)	Negative-binomial distribution	α RRA and MLE for the gene-level statistics
PBNPA (Jia et al. 2017)		Non-parametric permutation-test for each replicate
ScreenBeam (Yu et al. 2016)	Normal distribution	Bayesian hierarchical modeling
sgRSEA (Noh 2015)		Non-parametric permutation-test
PinAPL-Py (Spahn et al. 2017)	Negative-binomial model (control samples)	α RRA and STARS
RIGER (Luo et al. 2008)		Kolmogorov–Smirnov-based non-parametric statistics for the gene-level statistics
RSA (König et al. 2007)	hypergeometric distribution (sgRNA ranking)	Ranking-based statistics for the gene-level statistics