



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. 2019 29: 999-1008 originally 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)

References This article cites 47 articles, 11 of which can be accessed free at:
<http://genome.cshlp.org/content/29/6/999.full.html#ref-list-1>

Open Access Freely available online through the *Genome Research* Open Access option.

Creative Commons License This article, published in *Genome Research*, is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), 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](#).



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

Method

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

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

¹Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas 77030, USA; ²Jan and Dan Duncan Neurological Research Institute, Texas Children's Hospital, Houston, Texas 77030, USA; ³Department of Pediatrics, Baylor College of Medicine, Houston, Texas, USA; ⁴Howard Hughes Medical Institute, Houston, Texas 77030, USA

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

[Supplemental material is available for this article.]

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

The development and optimization of clustered regularly interspaced short palindromic repeats and CRISPR-associated protein 9 (CRISPR/Cas9) systems have propelled pooled screen approaches into even wider use. Besides the relative simplicity and low cost, the robustness of hit identification has reduced the requirement for redundancy in the number of targeting single-guide RNAs (sgRNAs), which allows the same size library to

be more diverse (Sanjana et al. 2014; Xu et al. 2015). Moreover, these pooled libraries have been made accessible through repositories such as Addgene (<https://www.addgene.org/>). CRISPR/Cas9 pooled screens are thus within the technical reach of most biomedical researchers (Doench 2017). For all this accessibility on the experimental side, however, the resulting bioinformatics data sets are enormous and complex. The analysis of thousands of genetic perturbations demands considerable bioinformatics expertise.

In our experience, most users are unaware of whether their tool of choice models the data according to Poisson, negative binomial, or Gaussian distribution, or of the relative strengths and weaknesses of these models. The current roster of tools bears the imprint of the history of RNA-seq data analysis: RNA-seq data were initially modeled using Poisson distributions (Marioni et al. 2008), which is a natural choice for simple read counts. Poisson distribution assumes that the mean and variance are equal, however, and with biological data, the variance is often greater than the mean. Analytic methods therefore turned to negative binomial distribution, which can handle overdispersed data (Anders and Huber 2010; Love et al. 2014). A number of popular tools still use negative binomial distribution to analyze sgRNA screen data (Li et al. 2014; Spahn et al. 2017), even though the structure of the sgRNA screen data is very different from that of RNA-seq. In the latter, there is huge variation in transcript lengths, from 60 bp to 2.4 Mbp, but all sgRNAs for any given gene are designed to have the same length. This often leads to the variance being less than the mean (Supplemental Fig. S1). We hypothesized that a

⁵Present address: University of Ottawa Brain and Mind Research Institute, Ottawa Institute of Systems Biology, and Department of Cellular and Molecular Medicine, University of Ottawa, Ottawa, Ontario K1H 8M5, Canada

Corresponding author: zhandong.liu@bcm.edu

Article published online before print. Article, supplemental material, and publication date are at <http://www.genome.org/cgi/doi/10.1101/gr.245571.118>. Freely available online through the *Genome Research* Open Access option.

© 2019 Jeong et al. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

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

We therefore developed CB², a new web-based algorithm that uses beta-binomial distribution, and compared its performance with that of the eight most commonly used algorithms (HiTSelect [Diaz et al. 2015], MAGeCK [Li et al. 2014], PBNPA [Jia et al. 2017], PinAPL-Py [Spahn et al. 2017], RIGER [Luo et al. 2009], RSA [König et al. 2007], ScreenBEAM [Yu et al. 2016], and sgrSEA [<https://cran.r-project.org/web/packages/sgrSEA/>]), which encompass both parametric and nonparametric approaches (Table 1). We applied all these methods to 10 different biological data sets, taken from fields ranging from cancer to basic cell biology (Koike-Yusa et al. 2014; Parnas et al. 2015; Evers et al. 2016; Golden et al. 2017; Li et al. 2018; Sanson et al. 2018).

Results

CB² is more sensitive in target gene identification than existing methods

Identifying candidates by a statistical hypothesis test is a key component in any screen analysis. In CB², we adapted a beta-binomial model (Baggerly et al. 2003) with a modified Student's *t*-test to measure differences in sgRNA levels, followed by Fisher's combined probability test (Fisher 1925) to estimate the gene-level significance. We chose Fisher's method for two reasons: first, to keep the entire pipeline parametric and, second, to keep CB² as fast as possible (robust rank aggregation [RRA] requires permuting the data, a nonparametric feature, so it runs slower than Fisher's method). Furthermore, when we compared Fisher's method against RRA, we found that RRA is not effective in combining the *P*-values estimated by beta-binomial distribution (Supplemental Fig. S2).

We compared CB² with eight state-of-the-art methods on three benchmark data sets evaluating gene essentiality (Evers et al. 2016) using different technologies: CRISPR nuclease gene knockout via Cas9 (CRISPRn) and CRISPRinterference (CRISPRi; a CRISPR/Cas9 system with a catalytically inactive Cas9 fused to the transcriptional repressor KRAB, which results in gene repression). These benchmark data sets (CRISPRn-RT112, CRISPRn-UMUC3, and CRISPRi-RT112) were constructed based on 46 genes that are essential for cell survival and 47 genes that are nonessential. We first tested whether these methods could easily distinguish between essential and nonessential genes. We found that each method clearly discriminates essentiality by their gene rankings (Supplemental Fig. S3A). In addition, gene rankings obtained

from each method, except PBNPA for the CRISPRi-RT112 data set, are highly correlated (R^2 is [0.86, 0.98] for CRISPRn-RT112, [0.85, 0.98] for CRISPRn-UMUC3, and [0.72, 0.96] for CRISPRi-RT112) (Supplemental Fig. S3B). CB², ScreenBEAM, and MAGeCK produced very similar gene rankings across all the benchmark data sets. We also compared the precision-recall (PR) and receiver operating characteristic (ROC) curves across the methods and calculated the area under the curve (AUC) of each. CB² recorded the best PR-AUC and ROC-AUC scores for both CRISPRn screen data sets (Supplemental Figs. S4, S5). HiTSelect had the best PR-AUC and ROC-AUC scores for the CRISPRi-RT112 data set, for which CB² achieved comparable scores (Supplemental Fig. S6). Although the gene ranking is similar among these methods (Supplemental Fig. S3), the estimated *P*-values and false-discovery rates (FDRs) are very different. These results highlight the importance of using FDR to guide the gene selection process.

Although several CRISPR screens (Zhou et al. 2014; Parnas et al. 2015; Aguirre et al. 2016) have prioritized candidate hits by ranking, they do not provide statistical estimates of error rates. These methods, therefore, rely on an arbitrary rule to select the top candidate genes and are prone to biased selections and high hit attrition rates. One solution is candidate selection by a quantitative statistical measure such as a *P*-value or a false-positive rate (FDR) cut-off. To assess the detection powers of FDR of established CRISPR screen analysis methods and CB², we measured the F1-score ($2 \times \text{precision} \times \text{recall} / (\text{precision} + \text{recall})$) of each method, namely, the harmonic average of the precision and the recall, for FDR thresholds ranging from 10% to 0.01%. CB² outperformed all other methods at every FDR cut-off level, and all other methods lost their detection powers at more rigorous FDRs (Fig. 1A,B; Supplemental Fig. S7). In other words, all methods showed a small type-I error owing to the strong lethality phenotype of the CRISPR assay, but CB² showed a significantly lower type-II error than the other methods (Supplemental Fig. S8). Across all paradigms tested with different FDR cut-offs, CB² performed the best, with a much larger F1-score and recall. Thus, CB² is both sensitive and specific in selecting candidate genes.

To understand the differences produced by these methods, we next tested a prototypical essential gene, *RPL5*, to compare the gene-level enrichment across data sets and analytical tools. In the first CRISPR screen on an RT112 cell line, we expected to see the depletion of sgRNAs targeting *RPL5* in group *T*₁. Out of the 10 sgRNAs that target this gene in the CRISPRn-RT112 data set, six showed a strong decrease in the group *T*₁. CB² estimated an FDR of 2.07×10^{-19} , whereas only three other methods (HiTSelect, ScreenBEAM, and sgrSEA) estimated FDR to be <0.01

Table 1. Statistical models used by CB² and existing methods

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 multiobjective 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)		Nonparametric permutation test for each replicate
ScreenBeam (Yu et al. 2016)	Normal distribution	Bayesian hierarchical modeling
sgrSEA ^a		Nonparametric 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 nonparametric statistics for the gene-level statistics
RSA (König et al. 2007)	Hypergeometric distribution (sgRNA ranking)	Ranking-based statistics for the gene-level statistics

All of the methods were used in the target identification benchmarking.

^a<https://cran.r-project.org/web/packages/sgrSEA/>.

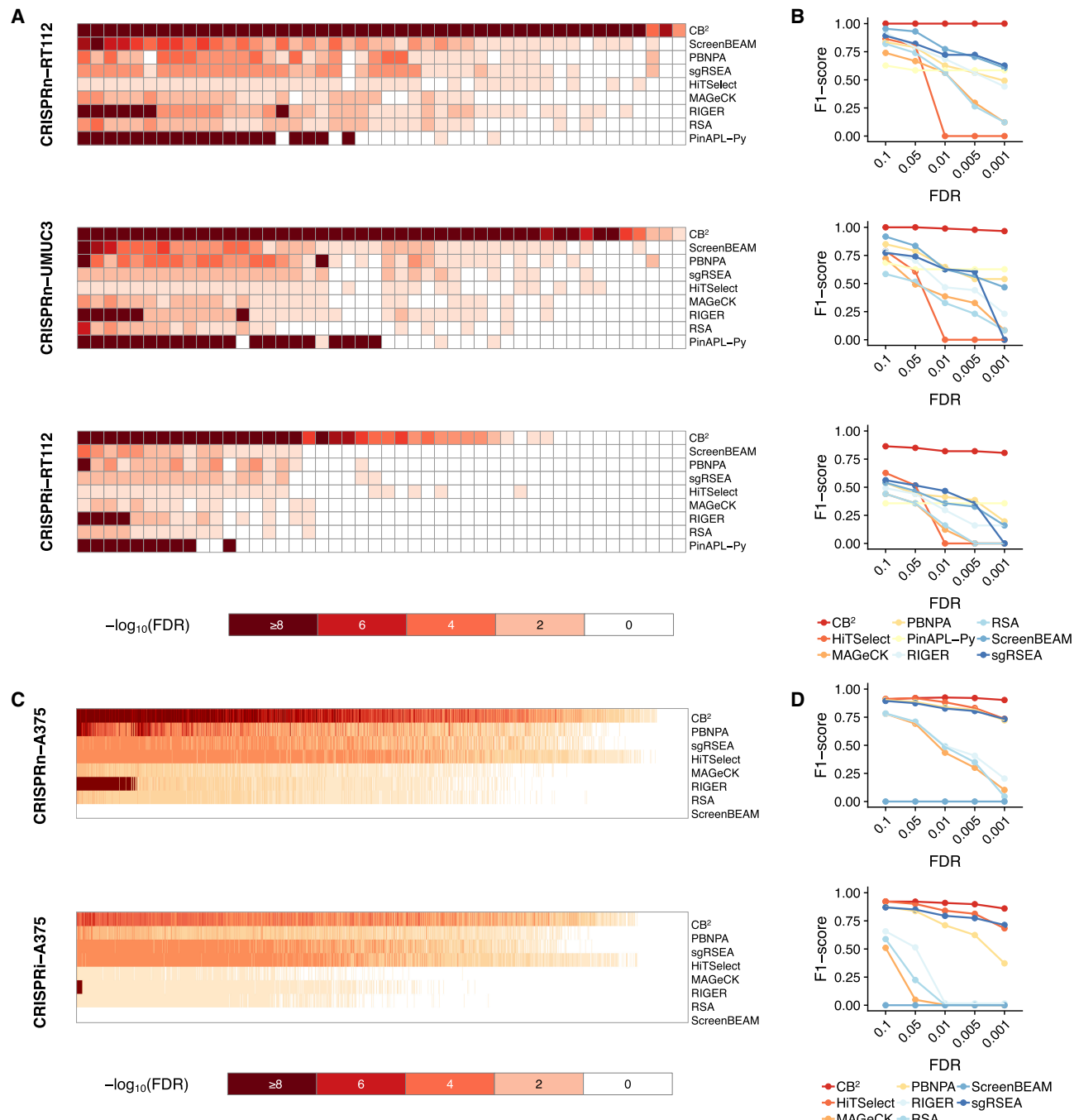


Figure 1. CB² offers robust target identification with high precision and recall. (A,B) Benchmark results using data from Evers et al. (2016). (A) Heatmaps illustrate FDRs of gene statistics from each of nine leading high-complexity pooled screen analysis tools. (B) F1-score measurements at different FDR cut-offs across all methods. At commonly used FDR cut-offs, CB² can identify most of the essential genes with high rates of precision and recall. (C,D) Same representation as in A and B, using data from Sanson et al. (2018).

(Fig. 2A). Next, we looked at the same gene in the UMUC3 cell line. Five of the 10 sgRNAs targeting *RPL5* decreased in group T₁, and all five sgRNAs were among those identified in RT112 cell line. CB² estimated an FDR of 3.81×10^{-10} for *RPL5*, whereas none of the other methods identified it to be significantly depleted with an FDR cut-off below 0.01 (Fig. 2B). Lastly, in the CRISPRi-RT112 data set, three of seven sgRNAs indicated depletions, but only CB² was able to estimate the FDR of 2.78×10^{-8} , and the other methods

did not count *RPL5* as a hit in the data set (Fig. 2C). Overall, CB² produced more reliable hit identification than other methods based on statistical cut-offs for the gene tested (*COPS8* and *RPL27*) (Supplemental Figs. S12, S13).

We performed the same analysis on two distinct data sets (Sanson et al. 2018) to determine how CB² performs compared with other methods for genome-wide screening analysis. Sanson et al. (2018) used novel optimized libraries for genome-wide

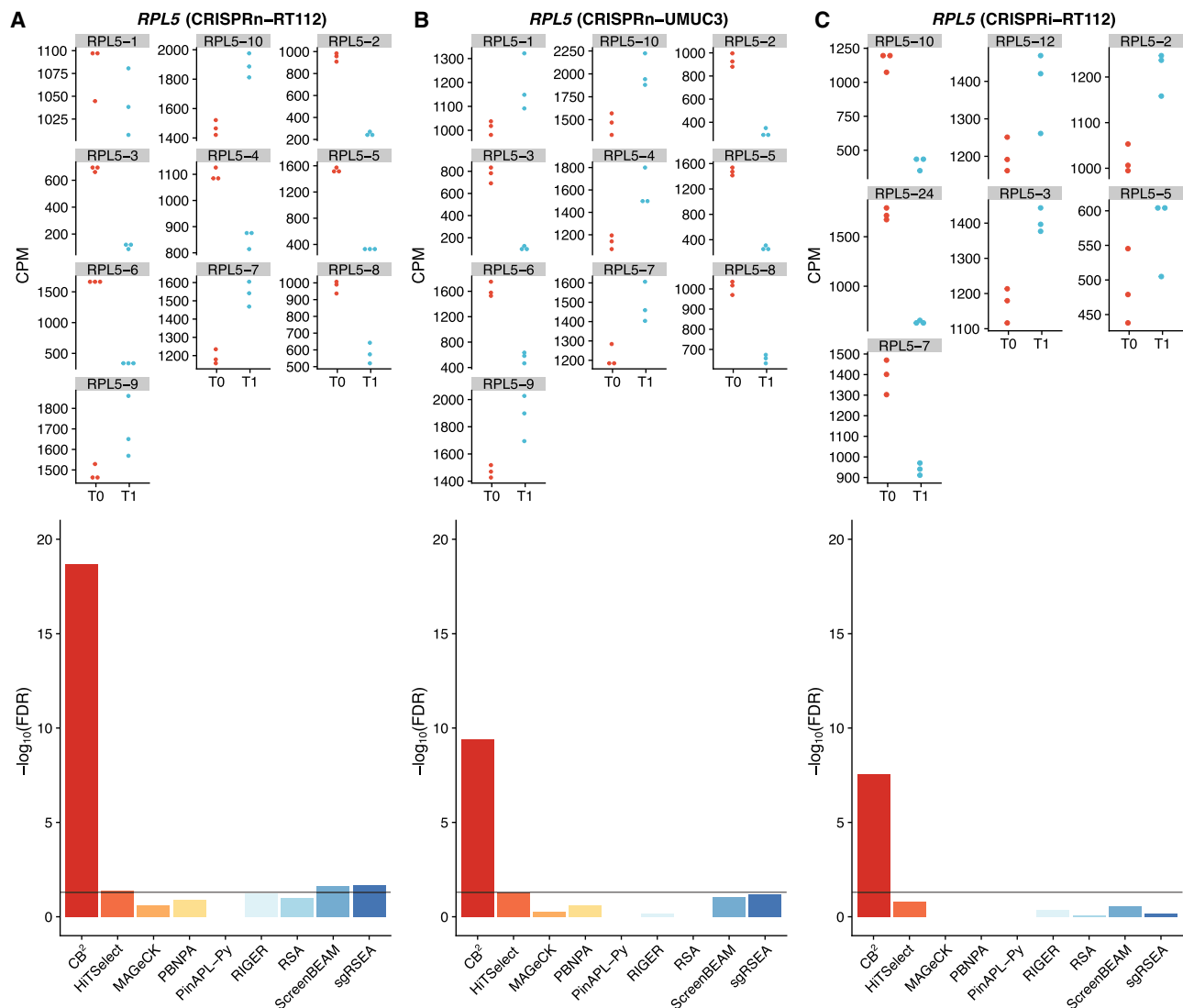


Figure 2. CB² detects essential genes missed by other leading methods: the case of *RPL5*. sgRNA quantification for *RPL5* in cell line (A) RT112, (B) UMUC3 using CRISPRn, and (C) RT112 using the CRISPRi library. The top panels show counts per million (CPM) of sgRNAs that target *RPL5* for each group (T_0 and T_1), and the bottom panels indicate the reported the FDR for *RPL5* in each screen across all the methods. A horizontal line at FDR=0.01 is used as a threshold for statistical cutoff. CB² outperforms all other methods of identifying *RPL5* as an essential gene across all benchmark data sets.

CRISPRn (Brunello), CRISPRi (Dolcetto), and CRISPRa (Calabrese) screening and showed these libraries outperform other previously established libraries, such as GeCKO (Sanjana et al. 2014) and hCRISPRi-v2 (Horlbeck et al. 2016). For our analysis, we chose two screens for benchmarking, both data sets from a dropout screen in A375 melanoma cells: One used the Brunello CRISPRn library with tracr-v2 tracrRNA (CRISPRn-A375); the other used the Dolcetto CRISPRi Set A library (CRISPRi-A375). Each data set contains a control sample (plasmid DNA) and three biological replicates. We used the gold-standard gene sets of 1580 essential and 927 nonessential genes reported by Hart et al. (2014, 2015) to assess the performance of the methods. (We excluded PinAPL-Py from benchmarking because it does not report statistics when the input contains only one control sample.) CB² outperformed other methods at the stringent FDR cutoff level (Fig. 2; Supplemental Fig. S9). CB² outperformed all other methods in F-1 and PR measures at the stringent FDR cut-offs on A375 genome-wide screen data sets of

Sanson et al. (2018). F1-score (top), precision (middle), and recall (bottom) for each method on two benchmark data sets are presented as a function of FDR cut-off values (Fig. 1C,D; Supplemental Fig. S9) and provided higher AUC values of PR and ROC curves than the other methods (Supplemental Figs. S10, S11).

These results indicate that CB² more accurately estimates the gene-level FDR. The use of FDR in selecting hits is critical in real data analysis because the arbitrary selection of top genes is purely heuristic. CB² is better at identifying true hits based on multiple concordant sgRNAs targeting the same gene.

CB² is more specific in target gene detection than existing methods

To test the idea that a beta-binomial model would better fit the data and more accurately identify changes in sgRNA, we compared the sgRNA level statistics on several CRISPR pooled libraries containing nontargeting sgRNAs as negative controls. Nontargeting

sgRNAs are not supposed to show any differential enrichment and can be used to assess the quality of the method. Parnas et al. (2015) used the Mouse CRISPR Knockout Pooled Library for their genome-wide screen (GeCKO v2; Addgene 100000052, 100000053), which contains 1000 nontargeting sgRNAs. We therefore used this data set to measure the specificity with which CB² and MAGeCK detect true negatives. We compared the unadjusted *P*-values for sgRNAs because the FDR is controlled at the gene level.

CB² showed greater specificity (the proportion of actual negatives that are correctly identified) than MAGeCK across a wide range of *P*-value thresholds. At a *P*-value threshold of 0.01, CB² had a specificity of 86% versus MAGeCK's 68% (Fig. 3A). Next, we plotted the log-fold change against the *P*-value levels in a volcano plot. The majority of the negative control sgRNAs were correctly detected by CB², whereas MAGeCK showed a one-side long tail for positively changed sgRNAs, producing inflated *P*-values for a group of negative controls (Fig. 3B). Many of these false positives showed extremely low *P*-values (ranging from 10⁻⁵ to 10⁻⁴⁰), indicating a strong selection bias. To understand the impact of this selection bias, we analyzed the rest of the sgRNA library with both methods. At the same threshold (*P* < 0.01), CB² selected 12,648 sgRNAs, whereas MAGeCK selected 31,381 sgRNAs; 2971 sgRNAs were identified by both methods (Fig. 4B). We applied the same analysis to the CRISPRn-A375 data set (Sanson et al. 2018), which contains 1000 nontargeting sgRNAs. CB² shows higher specificity than MAGeCK, except when setting a *P*-value cut-off at 0.2. Similar *P*-value distributions shown in Figure 3 were also found for this data set (Supplemental Fig. S14).

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

Next, we performed the same sgRNA-level comparisons on two additional data sets. In the first study, a differentiation screen was conducted to identify target genes that maintain naive pluripotency (Li et al. 2018) using the Mouse Improved Genome-wide Knockout CRISPR Library v2 (Addgene 67988). The library contains 91,319 sgRNAs targeting 18,542 mouse genes. To identify

differentially enriched sgRNAs, we kept the same threshold (*P* < 0.01) for both CB² and MAGeCK. CB² identified 732 sgRNAs whereas MAGeCK identified 5105 sgRNAs (395 sgRNAs shared between the two) (Fig. 4C,D), and we observed the same trend as in screening data set of Parnas et al. (2015) screening data set (Fig. 4A,B). We found the same trend on the data sets of Evers et al. (2016) (Supplemental Fig. S15). Thus, CB²'s accurate sgRNA-level statistics are attributable to its use of the beta-binomial model.

CB² provides more accurate alignment without parameter tuning

Many CRISPR pooled screens use in-house scripts to quantify sgRNA abundance (Gilbert et al. 2014; Golden et al. 2017; Iorio et al. 2018; Li et al. 2018) or other alignment algorithms for RNA-seq (Sanjana et al. 2014; Hart et al. 2015; Parnas et al. 2015; DeJesus et al. 2016). These codes are often not shared publicly and are not easily reusable. Both MAGeCK and PinAPL-Py provide an integrated mapping function, but PinAPL-Py requires complex parameter tuning and MAGeCK samples only the first million reads to estimate the location of sgRNAs in FASTQ files. Furthermore, there is no systematic comparison of mapping accuracy in the literature, and users lack reliable guidelines for selecting mapping tools. We therefore introduced an adaptive hash-mapping algorithm into CB² and tested all three methods on six published data sets (Supplemental Table S1).

CB² showed consistently greater mapping accuracy than MAGeCK and PinAPL-Py (Fig. 5A). To understand why, we studied the reads that are mapped by CB² but not MAGeCK or PinAPL-Py in the CRISPRn-RT112 data set (Evers et al. 2016). MAGeCK mapped 64% of the reads compared with 75% by CB² (Fig. 5A). This is primarily because MAGeCK estimates the trimming windows using the first *N* reads from the input (*N* is 100,000 by default). There is no guarantee that these windows are optimal for the rest of the input files. If a sgRNA locates outside of the precomputed windows, MAGeCK will fail to detect it (Fig. 5B). PinAPL-Py does not precompute sgRNA locations based on a subset of reads but uses Cutadapt (Martin 2011) for flexible trimming followed by the Bowtie 2-based alignment (Langmead and Salzberg 2012). We found that PinAPL-Py failed to identify some of the

sgRNAs because reads fail to align owing to the incorrect trimming from Cutadapt even under several different tuning parameters (Fig. 5B). This is likely because of the frequent indels that occurred in the 5' adapter sequence region of the reads (see Fig. 5B): Usually the quad-nucleotide sequence "CACC," which is part of the U6 promoter, precedes the sgRNA sequence. In contrast, the "GTTT" sequence, which is the first 4 nucleotides (nt) of the sgRNA scaffold sequence, was present in all the reads. Given the fidelity of the GTTT sequence, the sequences mapped by CB² but missed by other algorithms are likely accurate and not false positives. CB² is currently the only CRISPR/Cas9 online screen analysis tool with parameter-free mapping and high accuracy in sgRNA quantification.

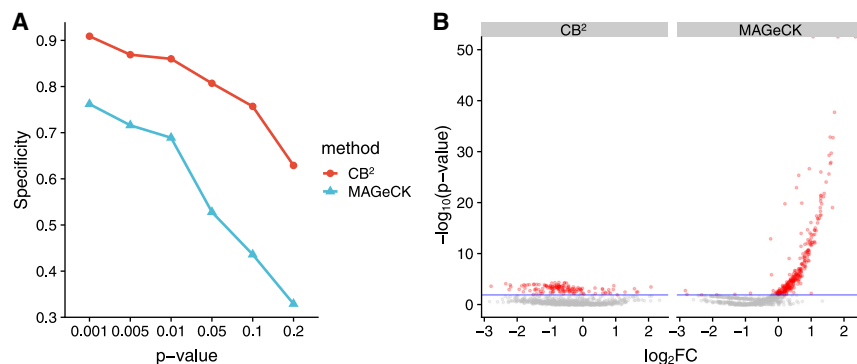


Figure 3. Comparison of false-positive rate for nontargeting sgRNAs on screen data of Parnas et al. (2015). (A) Specificity comparison between CB² and MAGeCK for the six different *P*-value thresholds. The y-axis indicates specificity, and the x-axis indicates the level of the *P*-value threshold for the specificity calculation. (B) Volcano plots of the *P*-value of nontargeting sgRNAs. The y-axis indicates the negative logarithm value of *P*-value, and the x-axis indicates the \log_2 value of fold-change. All of the data points are from negative control sgRNAs. False positives were plotted in red. Horizontal blue lines at *P* = 0.01 indicate the threshold for statistical cutoff.

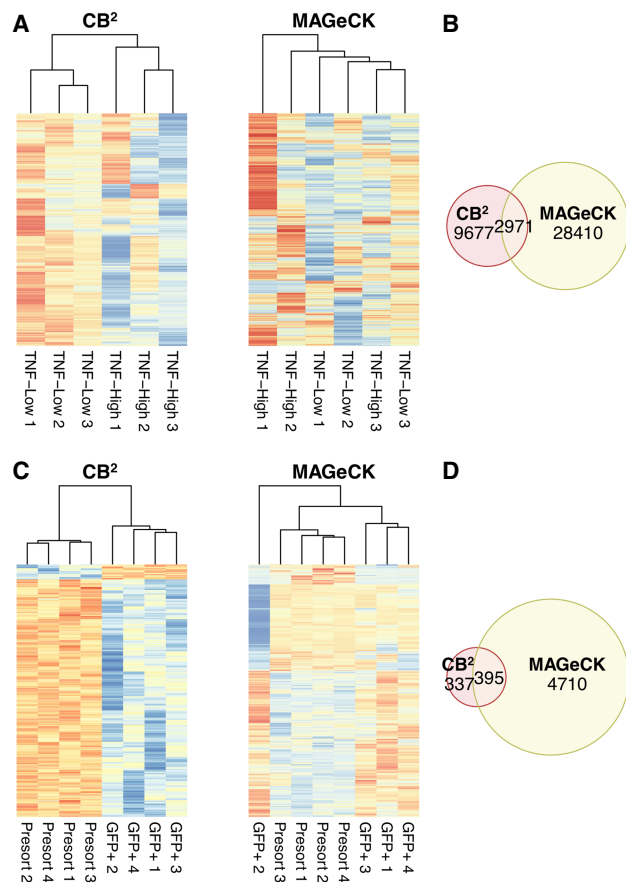


Figure 4. Discrepancy of sgRNA-level statistics between CB² and MAGeCK in two public CRISPR pooled screens. (A) Heatmaps of normalized read counts of the detected sgRNAs from the screen data of Parnas et al. (2015). (Left) A heatmap of sgRNAs detected by CB² only. (Right) A heatmap of sgRNAs detected by MAGeCK only. (B) Venn diagrams of sgRNAs detected by CB² and MAGeCK from the screen data of Parnas et al. (2015). (C, D) Same representations of A and B using data from Li et al. (2018).

CB² is accessible, secure, and easy to learn with CRISPRcloud

We had previously developed a web-based application called CRISPRcloud that could run any statistical testing and mapping algorithm through the cloud-based infrastructure provided by Amazon Web Service (AWS) (Jeong et al. 2017). We implemented CB² in the platform and added new features to increase speed and data security (Supplemental Fig. S16). CRISPRcloud is compatible with most modern web browsers (Google Chrome version 69 and later, Apple Safari version 11 and later, and Mozilla Firefox 48.0 and later) and operating systems (iOS, Windows, and Linux). Our fast client-side sgRNA mapping program reduces input files of several gigabytes into a single megabyte-sized file. By transferring a much smaller file through the Internet, CRISPRcloud decreases transfer time and prevents the sharing of raw input files, thereby eliminating downloading errors and data privacy issues in one step. Our adaptive mapping algorithm, via Angular (<https://angular.io/>) and TypeScript (<https://www.typescriptlang.org>), provides an open-source front-end web application platform. The enormous computing power needed to perform these operations mean that platforms built with a centralized server solution will have load-balancing problems when many users submit their requests

simultaneously, leading to the longer user waiting times and raising the risk of system-wide failure. CB² therefore provides a decentralized, cloud-computing-based, scalable service through a combination of AWS infrastructure that includes Amazon Elastic Compute Cloud (EC2) (<https://aws.amazon.com/ec2/>), Amazon Simple Storage Service (S3) (<https://aws.amazon.com/s3/>), and Amazon Simple Queue Service (SQS) (<https://aws.amazon.com/sqs/>). With this infrastructure, we launched a web service of CRISPRcloud. CRISPRcloud enables researchers with no programming background to preprocess, check the quality, statistically analyze, query, and visualize their CRISPR/Cas9 pooled screening data (Supplemental Table S2).

Discussion

The number of data sets for CRISPR/Cas9 screens in NCBI Gene Expression Omnibus have more than tripled each of the past 3 yr (39 data sets in 2015, 121 data sets in 2016, and 408 data sets in 2017). This expansion has outpaced the development of methods for analyzing the data, most of which use statistical models that are better suited to RNA-seq than to sgRNA data. Here we took into account the difference between the two types of data to develop a new algorithm, CB², and show that it is more sensitive, specific, and selective than eight other leading tools.

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

CB² also addressed another difficult task: quantifying sgRNA from NGS data. Except for the quantification algorithm provided by MAGeCK, most studies use in-house algorithms or extend established methods that were optimized for RNA-seq. CB² proved capable of fast and accurate alignment, with the ability to handle indels.

Last but not least is the challenge of making powerful tools readily accessible to the research community. Of the existing tools, PinAPL-Py (Spahn et al. 2017) and CRISPRcloud (Jeong et al. 2017) are the only two that support a graphical web interface and require no additional program installation. These programs are an important first step toward enabling the scientists who are actually generating the CRISPR/Cas9 screen data to analyze their large data sets. They still have limitations however: For PinAPL-Py, users still need to provide the adapter sequence to be trimmed, the error tolerance rate, and the quality threshold for trimmed reads. CRISPRcloud is the only framework into which the user can plug in any statistical tools or mapping algorithms, but

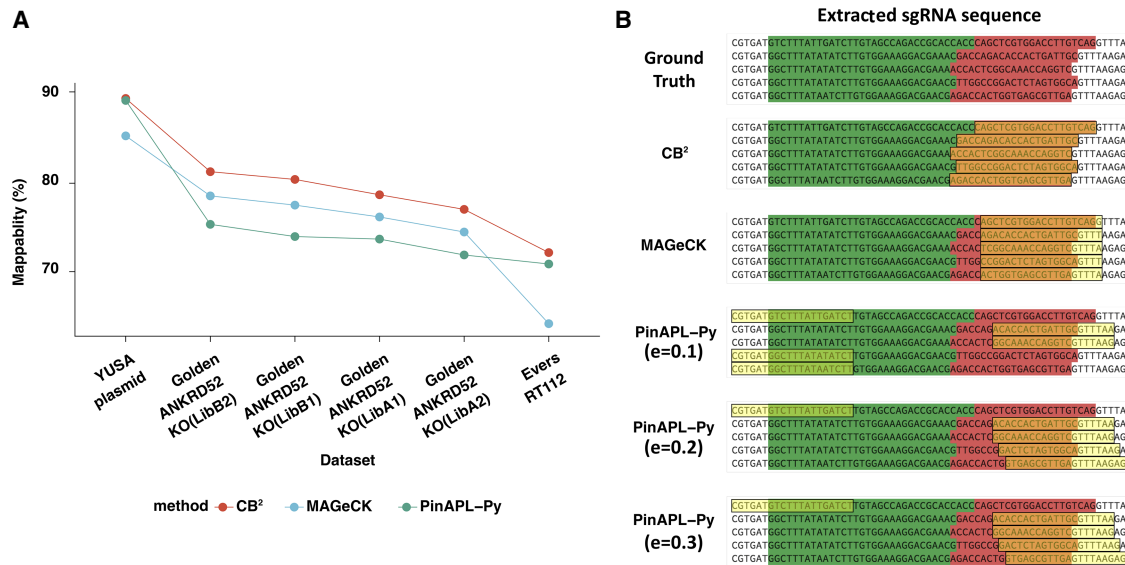


Figure 5. CB² outperforms MAGeCK and PinAPL-Py in the percentage of mapped reads over six benchmark data sets. (A) Read mappability of CB², MAGeCK, and PinAPL-Py across six different data sets. (B) Representative examples of reads that were not mapped by MAGeCK or PinAPL-Py. Adapters are highlighted with green; sgRNAs, with red. Yellow boxes show the predicted locations of sgRNAs by each method. Several parameters were used to optimize performances of PinAPL-Py.

transferring a large amount of sequence data over the Internet has inherent disadvantages such as long transfer times, vulnerability to file copying errors, and possible data security breaches. By taking advantage of the CRISPRcloud framework, CB² is fully web-based and designed to require only the minimal number of parameters for data analysis, because fewer parameters mean shorter learning curves for the majority of users. CB²'s power and accessibility will enable more laboratories to extract biologically relevant discoveries from CRISPR pooled screens.

Methods

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

We adapted a beta-binomial model proposed for serial analysis of gene expression (SAGE) by Baggerly et al. (2003). Specifically, let p_i be the true proportion of an sgRNA in sample i . We assume the value of p_i can vary from sample to sample and follows a beta distribution, $p_i \sim \text{Beta}(\alpha, \beta)$. Let X_i denote the number of read counts for a sgRNA in the i th sample. We assume X_i follows a binomial distribution, $X_i | p_i \sim \text{Binomial}(n_i, p_i)$, where n_i is the total number of mapped reads in sample i . To combine the estimated \hat{p}_i across multiple samples of the same treatment group, we proposed a linear 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. 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 w is estimated through gradient descent methods by minimizing the variance on p^A . Baggerly et al. (2003) showed that $w_i \propto [(1/(\alpha + \beta)) + (1/n_i)]^{-1}$.

CB² performs the sgRNA-level differential analysis between two groups using a Student's t -test-like statistic (Baggerly et al. 2003):

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

where p_A and p_B are the proportions of sgRNA, and V_A and V_B are the group variances of sgRNA, for groups A and B, respectively.

Test statistic t represents the strength of the difference of sgRNA abundance between groups A and B. In other words, a large positive t -value indicates that the quantity of sgRNA in group B is greater than in group A, and a large negative t -value indicates that the quantity of sgRNA in group B is less than that in group A.

The variance is estimated by

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

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

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

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

sgRNA p -value aggregation for gene-level statistics

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

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

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 the gene. χ^2 follows a chi-squared distribution with $2k$ degrees of freedom. To correct for multiple hypothesis testing, we adapted the Benjamini-Hochberg procedure (Benjamini and Hochberg 1995) to estimate the FDR.

Gene-level statistics benchmarking on existing methods

We used three different CRISPRn/CRISPRi pooled screen data sets (RT112 and UMUC3 cell line screens with CRISPR; RT112 cell line screen with CRISPRi) (Evers et al. 2016) and two different genome-wide CRISPRn/CRISPRi pooled screen data sets (A375 cell line screens) (Sanson et al. 2018), which provide ground-truth labels of essentiality for each gene. With those screening data sets and labels, we benchmarked the accuracy of essential gene detection by CB² with eight other published methods (HiTSelect, MAGeCK, PBNPA, PinAPL-Py, RIGER, RSA, ScreenBEAM, and sgrSEA). We computed the FDR for each gene from each method in the benchmark and set five different levels of FDR cut-off (0.1, 0.05, 0.01, 0.005, 0.001) for essential gene classification. For example, if we 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 below the cut-off value. We calculated recall (a recall value close to one indicates a prediction with a low false-negative rate), precision (a precision value close to one indicates a prediction with a low false-positive rate), and F-measure (the harmonic mean of PR) of all the methods at each FDR level. To allow each method to archive its best performance, we tuned parameters for the five parameter-tunable methods: MAGeCK (permutation parameter for RRA test), PBNPA (no.sim parameter), RIGER (alpha parameter), sgrSEA (multiplier parameter), and ScreenBEAM (burnin parameter) on CRISPRn-RT112 data set. The F-measure was used as a measure for the parameter tuning. Most of the methods showed robust performance regardless of the varied parameters, except sgrSEA and RIGER (Supplemental Fig. S17). Therefore, we used the default parameter for MAGeCK, PBNPA, and ScreenBEAM for other data sets and used an optimized parameter for sgrSEA and RIGER. We also calculated the AUC of the PR curves and ROC curves of all the methods with FDR values.

CB²

The CB² R package was used in the benchmarking (R Core Team 2019). Benchmarking of CB² was performed without parameter tuning because CB² is parameter free. FDR values for negative changes between two different groups from CB² statistical analysis were used for benchmarking.

HiTSelect

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

MAGeCK

MAGeCK version 0.5.8 was used for benchmarking. We ran MAGeCK with the “mageck test” command with the following parameters: “--norm-method,” “median,” and “--adjust-method” “fdr.” We performed 100 permutations for the modified robust ranking aggregation (α -RRA) algorithm to estimate the gene-level statistics on the benchmark data sets.

ScreenBEAM

The ScreenBEAM R package (version 1.0.0, <https://github.com/jyyu/ScreenBEAM>) was used for benchmarking “data.type” parameter was set as “NGS,” and “do.normalization” was set as TRUE, and “nitt” and “burnin” parameters for Bayesian computing were set at 15,000 and 5000. ScreenBEAM does not provide the one-sided p -value for negative selection, so for the FDR comparison with other methods, we changed the FDR of a gene to one if the β of the gene is greater than zero.

PinAPL-Py

We used the PinAPL-Py website (<http://pinapl-py.ucsd.edu>) to perform the benchmarking. For the sgRNA read counting, we used “GGCTTTATATATCTTGTGGAAAGGACGAAACACCG, GCTTTATATATCTTGTGGAAAGGACGAAACACCG,” and “CTTTATATATCTTGTGGAAAGGACGAAACACCG” for “seq_5_end” parameters of “CRISPRn-RT112,” “CRISPRn-UMUC3,” and “CRISPRi-RT112” data sets. We used CPM normalization and set the GeneMetric parameter as “aRRA” to perform a modified robust ranking aggregation (α -RRA). We used the combined FDR values for each gene in the benchmarking.

RIGER

We used the Java implementation of RIGER (version 2.0; <https://github.com/broadinstitute/rigerj>) to perform the benchmarking. We set the “alpha” parameter at 0.1 on the Evers et al. (2016) data sets and at 1.0 on the Sanson et al. (2018) data sets. log₂ fold-change values calculated by CB² were used as an input of RIGER.

RSA

We used the Python implementation of RSA (version 1.9, <https://admin-ext.gmf.org/publications/RSA/>). log₂ fold-change values calculated by CB² were used as an input of RSA.

sgRSEA

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

PBNPA

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

Specificity measure at sgRNA level

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

Algorithm for quantifying sgRNA abundance

Previous sgRNA abundance quantification methods

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

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

However, as pointed out by Spahn et al. (2017), CRISPRcloud quantification algorithm can produce erroneous mapping results if the sgRNA sequences are staggered, because CRISPRcloud

extracts sgRNA sequence for each read at a fixed location. Another limitation of CRISPRcloud is the fact that the user must decide where the extraction site is. Nevertheless, CRISPRcloud did not require tuning and was thus arguably more user-friendly than other tools. For instance, in PinAPL-Py, users need to set many tuning parameters for sgRNA quantification: adapter error rate parameter for trimming, matching and ambiguity thresholds, and parameters for alignment seed for the Bowtie alignment (Spahn et al. 2017).

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

The binary representation of sgRNA sequence lowers the cost of computation

We used a binary representation for sgRNA sequence. This approach is memory efficient and improves the user experience at the client side (Melsted and Pritchard 2011). It only needs $max(K, 2M)$ bits to store an sgRNA-sequence, where M is the length of the sequence, and K is the bit size to store a primitive integer in the machine (usually 64 bits) because we only need two bits to save a nucleotide (i.e., “A” is “00,” “C” is “01,” “G” is “11,” and “T” is “10”). The memory size is about half of that required for storing a character string of the sequence; that is, 160 bits are needed to store a 20-nt sgRNA sequence. Another benefit of binary representation is that it lowers the time complexity for the shift operator when comparing all k -mers of an sgRNA read using a sliding window. This is an essential function for the quantification algorithm in CB². Compared with the string shift operator functions, such as string copy, substring extraction, and concatenation, the binary representation produces significantly shorter running times. Algorithm 2 in Supplemental Methods illustrates how a sgRNA library converted to a bit sequence and stored the converted sequence into a hash table.

Sliding window–based algorithm gives a high-resolution quantification with comparable running time

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

For the case of a reverse complement sequenced sample, the entire procedure is repeated on the reverse complement reference sgRNA library, and the reads are scanned from 3′ to 5′. After both assays are performed (5′ to 3′ and 3′ to 5′ with the reverse complement reference sequences), mapping results between both sequences are compared. The one with a larger count corresponds to the correct sequence mapping. We compared the mappability of CB² to those of MAGeCK (Li et al. 2015) and PinAPL-Py (Spahn et al. 2017) across multiple data sets from previous studies (Fig. 5).

Software availability

CRISPRBetaBinomial or CB² is available as Supplemental Code and at <https://CRAN.R-project.org/package=CB2>. All of the data and scripts for the benchmarking are available in the Supplemental Material and at GitHub (<https://github.com/hyunhwaj/CB2-Experiments>). Parameters used in these experiments are described above in the Methods. CRISPRcloud is available at <http://crispr.nrihub.org>.

Acknowledgments

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

Author contributions: 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 software. H-H.J. and Z.L. performed analysis. Z.L. supervised the project. H-H.J., M.W.C.R., and Z.L. wrote the manuscript with input from all the authors.

References

- Aguirre AJ, Meyers RM, Weir BA, Vazquez F, Zhang C-Z, Ben-David U, Cook A, Ha G, Harrington WF, Doshi MB, et al. 2016. Genomic copy number dictates a gene-independent cell response to CRISPR/Cas9 targeting. *Cancer Discov* **6**: 914–929. doi:10.1158/2159-8290.CD-16-0154
- Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol* **11**: R106. doi:10.1186/gb-2010-11-10-r106
- Baggerly KA, Deng L, Morris JS, Aldaz CM. 2003. Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics* **19**: 1477–1483. doi:10.1093/bioinformatics/btg173
- Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Ser B Methodol* **57**: 289–300. doi:10.2307/2346101
- DeJesus R, Moretti F, McAllister G, Wang Z, Bergman P, Liu S, Frias E, Alford J, Reece-Hoyes JS, Lindeman A, et al. 2016. Functional CRISPR screening identifies the ufm1ylation pathway as a regulator of SQSTM1/p62. *eLife* **5**: e17290. doi:10.7554/eLife.17290
- Diaz AA, Qin H, Ramalho-Santos M, Song JS. 2015. HiTSelect: a comprehensive tool for high-complexity-pooled screen analysis. *Nucleic Acids Res* **43**: e16. doi:10.1093/nar/gku1197
- Doench JG. 2017. Am I ready for CRISPR? A user’s guide to genetic screens. *Nat Rev Genet* **19**: 67–80. doi:10.1038/nrg.2017.97
- Dutta B, Azhir A, Merino L-H, Guo Y, Revanur S, Madhamshettiar PB, Germain RN, Smith JA, Simpson KJ, Martin SE, et al. 2016. An interactive web-based application for comprehensive analysis of RNAi-screen data. *Nat Commun* **7**: 10578. doi:10.1038/ncomms10578
- Evers B, Jastrzebski K, Heijmans JPM, Grerum W, Beijersbergen RL, Bernards R. 2016. CRISPR knockout screening outperforms shRNA and CRISPRi in identifying essential genes. *Nat Biotechnol* **34**: 631–633. doi:10.1038/nbt.3536
- Fisher RA. 1925. *Statistical methods for research workers*. Oliver and Boyd, Edinburgh.
- Gilbert LA, Horlbeck MA, Adamson B, Villalta JE, Chen Y, Whitehead EH, Guimaraes C, Panning B, Ploegh HL, Bassik MC, et al. 2014. Genome-scale CRISPR-mediated control of gene repression and activation. *Cell* **159**: 647–661. doi:10.1016/j.cell.2014.09.029
- Golden RJ, Chen B, Li T, Braun J, Manjunath H, Chen X, Wu J, Schmid V, Chang T-C, Kopp F, et al. 2017. An Argonaute phosphorylation cycle promotes microRNA-mediated silencing. *Nature* **542**: 197–202. doi:10.1038/nature21025
- Hart T, Brown KR, Sircoulomb F, Rottapel R, Moffat J. 2014. Measuring error rates in genomic perturbation screens: gold standards for human functional genomics. *Mol Syst Biol* **10**: 733. doi:10.15252/msb.20145216
- Hart T, Chandrashekar M, Aregger M, Steinhardt J, Brown KR, MacLeod G, Mis M, Zimmermann M, Fradet-Turcotte A, Sun S, et al. 2015. High-

- resolution CRISPR screens reveal fitness genes and genotype-specific cancer liabilities. *Cell* **163**: 1515–1526. doi:10.1016/j.cell.2015.11.015
- Horlbeck MA, Gilbert LA, Villalta JE, Adamson B, Pak RA, Chen Y, Fields AP, Park CY, Corn JE, Kampmann M, et al. 2016. Compact and highly active next-generation libraries for CRISPR-mediated gene repression and activation. *eLife* **5**: e19760. doi:10.7554/eLife.19760
- Hu G, Luo J. 2012. A primer on using pooled shRNA libraries for functional genomic screens. *Acta Biochim Biophys Sin* **44**: 103–112. doi:10.1093/abbs/gmr116
- Iorio F, Behan FM, Gonçalves E, Bhosle SG, Chen E, Shepherd R, Beaver C, Ansari R, Pooley R, Wilkinson P, et al. 2018. Unsupervised correction of gene-independent cell responses to CRISPR-Cas9 targeting. *BMC Genomics* **19**: 604. doi:10.1186/s12864-018-4989-y
- Jeong HH, Kim SY, Rousseaux MW, Zoghbi HY, Liu Z. 2017. CRISPRcloud: a secure cloud-based pipeline for CRISPR pooled screen deconvolution. *Bioinformatics* **33**: 2963–2965. doi:10.1093/bioinformatics/btx335
- Jia G, Wang X, Xiao G. 2017. A permutation-based non-parametric analysis of CRISPR screen data. *BMC Genomics* **18**: 545. doi:10.1186/s12864-017-3938-5
- Koike-Yusa H, Li Y, Tan E-P, Velasco-Herrera MDC, Yusa K. 2014. Genome-wide recessive genetic screening in mammalian cells with a lentiviral CRISPR-guide RNA library. *Nat Biotechnol* **32**: 267–273. doi:10.1038/nbt.2800
- König R, Chiang C, Tu BP, Yan SF, DeJesus PD, Romero A, Bergauer T, Orth A, Krueger U, Zhou Y, et al. 2007. A probability-based approach for the analysis of large-scale RNAi screens. *Nat Methods* **4**: 847–849. doi:10.1038/nmeth1089
- Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**: 357–359. doi:10.1038/nmeth.1923
- Li W, Xu H, Xiao T, Cong L, Love MI, Zhang F, Irizarry RA, Liu JS, Brown M, Liu XS. 2014. MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9 knockout screens. *Genome Biol* **15**: 554. doi:10.1186/s13059-014-0554-4
- Li W, Köster J, Xu H, Chen C-H, Xiao T, Liu JS, Brown M, Liu XS. 2015. Quality control, modeling, and visualization of CRISPR screens with MAGeCK-VISPR. *Genome Biol* **16**: 281. doi:10.1186/s13059-015-0843-6
- Li M, Yu JSL, Tilgner K, Ong SH, Koike-Yusa H, Yusa K. 2018. Genome-wide CRISPR-KO screen uncovers mTORC1-mediated Gsk3 regulation in naive pluripotency maintenance and dissolution. *Cell Rep* **24**: 489–502. doi:10.1016/j.celrep.2018.06.027
- Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**: 550. doi:10.1186/s13059-014-0550-8
- Luo B, Cheung HW, Subramanian A, Sharifnia T, Okamoto M, Yang X, Hinkle G, Boehm JS, Beroukhim R, Weir BA, et al. 2008. Highly parallel identification of essential genes in cancer cells. *Proc Natl Acad Sci* **105**: 20380–20385. doi:10.1073/pnas.0810485105
- Luo J, Emanuele MJ, Li D, Creighton CJ, Schlabach MR, Westbrook TF, Wong K-K, Elledge SJ. 2009. A Genome-wide RNAi screen identifies multiple synthetic lethal interactions with the Ras oncogene. *Cell* **137**: 835–848. doi:10.1016/j.cell.2009.05.006
- Marioni JC, Mason CE, Mane SM, Stephens M, Gilad Y. 2008. RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays. *Genome Res* **18**: 1509–1517. doi:10.1101/gr.079558.108
- Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* **17**: 10. doi:10.14806/ej.17.1.200
- Melsted P, Pritchard JK. 2011. Efficient counting of *k*-mers in DNA sequences using a bloom filter. *BMC Bioinformatics* **12**: 333. doi:10.1186/1471-2105-12-333
- Mohr SE, Smith JA, Shamu CE, Neumüller RA, Perrimon N. 2014. RNAi screening comes of age: improved techniques and complementary approaches. *Nat Rev Mol Cell Biol* **15**: 591–600. doi:10.1038/nrm3860
- Paddison PJ, Silva JM, Conklin DS, Schlabach M, Li M, Aruleba S, Balija V, O'Shaughnessy A, Gnoj L, Scobie K, et al. 2004. A resource for large-scale RNA-interference-based screens in mammals. *Nature* **428**: 427–431. doi:10.1038/nature02370
- Park J, Al-Ramahi I, Tan Q, Mollema N, Diaz-Garcia JR, Gallego-Flores T, Lu H-C, Lagalwar S, Duvick L, Kang H, et al. 2013. RAS–MAPK–MSK1 pathway modulates ataxin 1 protein levels and toxicity in SCA1. *Nature* **498**: 325–331. doi:10.1038/nature12204
- Parnas O, Jovanovic M, Eisenhaure TM, Herbst RH, Dixit A, Ye CJ, Przybylski D, Platt RJ, Tirosh I, Sanjana NE, et al. 2015. A genome-wide CRISPR screen in primary immune cells to dissect regulatory networks. *Cell* **162**: 675–686. doi:10.1016/j.cell.2015.06.059
- R Core Team. 2019. *R: a language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna. <https://www.R-project.org/>.
- Sanjana NE, Shalem O, Zhang F. 2014. Improved vectors and genome-wide libraries for CRISPR screening. *Nat Methods* **11**: 783–784. doi:10.1038/nmeth.3047
- Sanson KR, Hanna RE, Hegde M, Donovan KF, Strand C, Sullender ME, Vaimberg EW, Goodale A, Root DE, Piccioni F, et al. 2018. Optimized libraries for CRISPR-Cas9 genetic screens with multiple modalities. *Nat Commun* **9**: 5416. doi:10.1038/s41467-018-07901-8
- Schlabach MR, Luo J, Solimini NL, Hu G, Xu Q, Li MZ, Zhao Z, Smogorzewska A, Sowa ME, Ang XL, et al. 2008. Cancer proliferation gene discovery through functional genomics. *Science* **319**: 620–624. doi:10.1126/science.1149200
- Shalem O, Sanjana NE, Hartenian E, Shi X, Scott DA, Mikkelsen TS, Heckl D, Ebert BL, Root DE, Dönnch JG, et al. 2014. Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science* **343**: 84–87. doi:10.1126/science.1247005
- Shao DD, Tsherniak A, Gopal S, Weir BA, Tamayo P, Stransky N, Schumacher SE, Zack TI, Beroukhim R, Garraway LA, et al. 2013. ATARIS: computational quantification of gene suppression phenotypes from multisample RNAi screens. *Genome Res* **23**: 665–678. doi:10.1101/gr.143586.112
- Silva JM, Marran K, Parker JS, Silva J, Golding M, Schlabach MR, Elledge SJ, Hannon GJ, Chang K. 2008. Profiling essential genes in human mammary cells by multiplex RNAi screening. *Science* **319**: 617–620. doi:10.1126/science.1149185
- Simon MM, Moresco EMY, Bull KR, Kumar S, Mallon A-M, Beutler B, Potter PK. 2015. Current strategies for mutation detection in phenotype-driven screens utilising next generation sequencing. *Mamm Genome* **26**: 486–500. doi:10.1007/s00335-015-9603-x
- Spahn PN, Bath T, Weiss RJ, Kim J, Esko JD, Lewis NE, Harismendy O. 2017. PinAPL-Py: a comprehensive web-application for the analysis of CRISPR/Cas9 screens. *Sci Rep* **7**: 15854. doi:10.1038/s41598-017-16193-9
- Wang T, Wei JJ, Sabatini DM, Lander ES. 2014. Genetic screens in human cells using the CRISPR-Cas9 system. *Science* **343**: 80–84. doi:10.1126/science.1246981
- Winter J, Schwering M, Pelz O, Rauscher B, Zhan T, Heigwer F, Boutros M. 2017. CRISPRAnalyzeR: interactive analysis, annotation and documentation of pooled CRISPR screens. bioRxiv doi:10.1101/109967
- Xu H, Xiao T, Chen CH, Li W, Meyer CA, Wu Q, Wu D, Cong L, Zhang F, Liu JS, et al. 2015. Sequence determinants of improved CRISPR sgRNA design. *Genome Res* **25**: 1147–1157. doi:10.1101/gr.191452.115
- Yu J, Silva J, Califano A. 2016. ScreenBEAM: a novel meta-analysis algorithm for functional genomics screens via Bayesian hierarchical modeling. *Bioinformatics* **32**: 260–267. doi:10.1093/bioinformatics/btv556
- Zhou Y, Zhu S, Cai C, Yuan P, Li C, Huang Y, Wei W. 2014. High-throughput screening of a CRISPR/Cas9 library for functional genomics in human cells. *Nature* **509**: 487–491. doi:10.1038/nature13166

Received October 23, 2018; accepted in revised form April 3, 2019.