



OPEN

Uncovering hidden cancer self-dependencies through analysis of shRNA-level dependency scores

Zohreh Toghrayee^{1,2} & Hesam Montazeri¹✉

Large-scale short hairpin RNA (shRNA) screens on well-characterized human cancer cell lines have been widely used to identify novel cancer dependencies. However, the off-target effects of shRNA reagents pose a significant challenge in the analysis of these screens. To mitigate these off-target effects, various approaches have been proposed that aggregate different shRNA viability scores targeting a gene into a single gene-level viability score. Most computational methods for discovering cancer dependencies rely on these gene-level scores. In this paper, we propose a computational method, named NBDep, to find cancer self-dependencies by directly analyzing shRNA-level dependency scores instead of gene-level scores. The NBDep algorithm begins by removing known batch effects of the shRNAs and selecting a subset of concordant shRNAs for each gene. It then uses negative binomial random effects models to statistically assess the dependency between genetic alterations and the viabilities of cell lines by incorporating all shRNA dependency scores of each gene into the model. We applied NBDep to the shRNA dependency scores available at Project DRIVE, which covers 26 different types of cancer. The proposed method identified more well-known and putative cancer genes compared to alternative gene-level approaches in pan-cancer and cancer-specific analyses. Additionally, we demonstrated that NBDep controls type-I error and outperforms statistical tests based on gene-level scores in simulation studies.

Large-scale CRISPR-Cas9 and RNAi screens have been increasingly used in cancer research to identify novel cancer vulnerabilities and therapeutic choices. While CRISPR-Cas9 technology can be used to perform knockout of gene function at the DNA level through multiple single-guide RNAs, RNAi screens knockdown genes at the mRNA level using a pool of short hairpin RNAs (shRNA)^{1–6}. Despite the potential benefits of these screens in identifying cancer dependencies, a major challenge in the analysis of pooled shRNA screens is to account for the off-target effects of reagents^{7,8}. Various computational approaches have been proposed to mitigate the impact of off-target effects in the analysis of shRNA screens^{9–12}. A common theme among all the previous approaches for identifying cancer dependencies is to first aggregate individual reagent effects into a single gene-level score using various computational tools such as RSA¹², ATARIS⁹, and DEMETER¹³. Gene-level scores are then used to compare subjects with highly diverse molecular profiles in order to infer cancer dependencies.

The RSA method employs a probabilistic approach to calculate absolute gene-level viability scores from multiple siRNAs targeting a specific gene. It evaluates whether the siRNAs targeting the gene are unusually top-ranked among all siRNAs in the screen¹². The ATARIS method provides relative gene-level dependency values by only incorporating a subset of RNAi reagents whose phenotypic effects are concordant across multiple samples⁹. DEMETER is another computational framework for estimating relative gene-level scores using multiple shRNA effects, assuming that each observed shRNA value is a linear combination of the corresponding gene-level effect and the batch effect of the corresponding seed sequence¹³. The gene-level effects are estimated using a stochastic gradient descent algorithm to minimize a regularized objective function. DEMETER2 extends the original DEMETER by using a hierarchical model for the gene and seed effects that integrate information across cell lines. DEMETER2 provides absolute gene-level scores and uses a Bayesian inference method for the parameter estimation¹¹. The gespeR method uses a regression model to account for sequence-dependent off-target effects, based on the TargetScan model for predicting relationships between siRNA and its off-targets¹⁰. TargetScan is a miRNA target prediction model that predicts mRNA fold change between wild-type and knockout cells based on various features of miRNA sequence¹⁴.

¹Department of Bioinformatics, Institute Biochemistry and Biophysics, University of Tehran, Tehran, Iran. ²Department of Bioinformatics, Kish International Campus University of Tehran, Kish, Iran. ✉email: hesam.montazeri@ut.ac.ir

Discovering cancer dependencies, including synthetic lethality and self-dependencies, is of crucial significance for identifying new cancer treatment options. Synthetic lethality refers to the interaction of two genes, where the simultaneous loss of function through either genetic events or inhibition results in cell death, but the loss of function of either gene alone does not. Several computational methods have been previously proposed for identifying synthetic lethality interactions using loss-of-function RNAi and CRISPR screens^{15–19}. Cancer self-dependency refers to a dependency type in which the loss of function of a specific gene leads to cell death preferentially in cells with specific molecular characteristics, such as mutations in the same gene. Studies have previously examined self-dependency in relation to missense and damaging mutations, as well as copy-number amplification, and have identified novel putative cancer genes through these investigations^{9,20,21}.

To identify cancer dependency through gene perturbation, several approaches first aggregate scores at the gene level and then perform statistical tests at this level^{9,11,12,22,23}. To the best of our knowledge, in analyzing RNAi screening data, the method proposed in this study is the only one that directly performs statistical tests at the shRNA level to identify cancer dependencies. However, for CRISPR analysis, MAGeCK method performs hypothetical testing at the sgRNA level²⁴. Similarly, in protein expression data analysis, it has been shown that analyzing data at the peptide level offers greater sensitivity and specificity^{25–27}. Additionally, there are methods tested at the shRNA level, such as siMEM, which is a time series-based design in RNAi screening, and Screen-Beam, a meta-analysis method that uses both shRNAs and sgRNAs. However, these methods are not comparable to this research^{28,29}.

In this research, we investigated the potential of using shRNA-level viability scores instead of gene-level data to enhance the statistical power for identifying cancer self-dependencies. To tackle the challenge of off-target effects, we developed a statistical method for analyzing dependency scores from perturbation screens at the shRNA level. Our hypothesis was that analyzing at the shRNA level would result in greater statistical power for uncovering hidden cancer dependencies. To evaluate the performance of our proposed method, we analyzed pan-cancer and cancer-specific analyses on the Project DRIVE dataset, an RNAi screening project that used deep coverage shRNA lentiviral libraries to target genes across 398 cell lines, providing a valuable resource for exploring and discovering cancer dependencies. In addition, we conducted simulation studies to evaluate the type-I error and statistical power of our approach across various sample sizes.

Materials and methods

shRNA viability data from Project DRIVE

In this research, we used the Project DRIVE data to find novel cancer dependencies. Project DRIVE conducted knockdown experiments in three different pools namely poolA, poolB, and BGPD using 158,114 shRNAs on 9850 genes in 398 cell lines across 39 cancer types by using a median of 20 pooled shRNA per gene. Pools poolA, poolB, and BGPD included 3492, 3577, and 4178 genes, respectively (Supplementary Fig. S1b). The shRNA viability scores were defined using next-generation sequencing counts as log fold change of shRNA read counts 14 days after the onset of the screen compared to shRNA initial abundance in the input library. For cancer-specific analyses, we only considered 26 cancer types that included at least four cell lines, following completion of all preprocessing steps.

Gene-level viability data from Project DRIVE

We used ATARiS and DEMETER2 methods for aggregating shRNA viability scores into gene-level dependency scores. ATARiS uses a subset of shRNAs with a consistent pattern of viability scores across cell lines and provides gene-level dependency scores relative to screened cell lines due to using median-centered shRNA viability scores. The ATARiS scores of Project DRIVE are available for 6557 genes in 398 cell lines. On the contrary, DEMETER2 scores are absolute dependency scores computed using a hierarchical Bayesian inference scheme by explicitly modeling the off-target effects associated with the seed sequence of each shRNA. The DEMETER2 scores of Project DRIVE are available for 7975 genes in 397 cell lines.

Copy number and mutation data

We downloaded the mutation data from the DepMap website and the copy number GISTIC2 data from the CBioPortal website³⁰. We considered three different classes of genetic alterations: non-missense mutations, missense mutations, and copy-number amplification. A gene has a non-missense mutation in a cell line if it harbors any of start codon deletion, stop codon deletion, start codon insertion, start codon insertion, frameshift deletion, frameshift insertion, in-frame deletion, in-frame insertion, nonsense mutation, and splice site. Missense mutations are simply defined as mutations annotated with missense mutation. We used the GISTIC2 to determine the copy-number status of a gene. We specifically used GISTIC2 values 2 and – 2 representing copy number amplification and deep deletion, respectively. Finally, a gene is considered wild-type in a cell line if it does not harbor missense, non-missense, and copy-number amplification and deep deletion.

TargetScan data

We used thermodynamic stability of seed pairing in shRNAs in RNAi screening to incorporate the off-target effects of seed sequences. The thermodynamic stability values of 7-mer seeds of the DRIVE data were extracted from TargetScan³¹. The total number of seeds in the TargetScan data and Project DRIVE are 16,384 (4⁷) and 13,071, respectively.

IntOGen data

We used the IntOGen data (version 2020.02.01) for validation of our method. It consists of 568 genes involved in cancer among which 422 and 474 genes were also available in the ATARiS and DEMETER2 scores, respectively³².

NBDep algorithm

In this section, we explain our proposed framework for identification of cancer self-dependencies, i.e. observing reduced viabilities preferentially in mutated cell lines in the respective gene upon knocking down (Fig. 1). In particular, we identified genes showing dependencies linked with three classes of genetic alterations namely missense mutation, non-missense mutation, copy number amplification as previously studied in Ref.²¹. The NBDep algorithm is as follows.

Step 1: preprocessing

We used the following preprocessing steps to select a subset of shRNAs, genes, and cell lines for further analysis:

1. We eliminated hypermutated cell lines, which were defined as cell lines with a mutation burden three standard deviations above the mean mutation burden across all cell lines.
2. We excluded genes without mutation or copy number data.
3. We eliminated shRNAs from the analysis if the initial abundance level was missing or below 50.
4. In some cases, multiple read counts were reported for a shRNA in a cell line. In these cases, we only retained the average abundance value for the cell line.

The above preprocessing steps resulted in data on 139,407 shRNAs targeting 7324 genes in 339 cell lines.

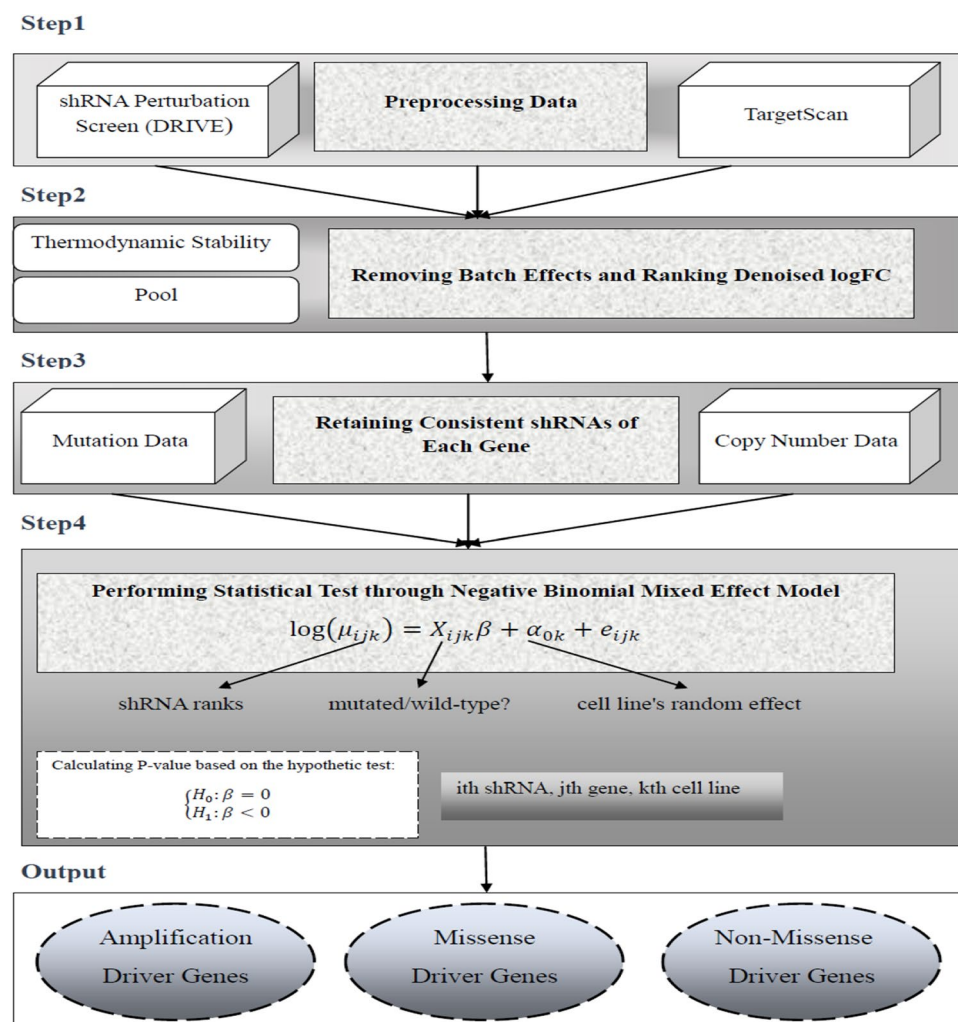


Figure 1. The flowchart of the NBDep algorithm. NBDep employs Project DRIVE, CCLE, GISTIC, and thermodynamic stability of shRNAs seeds obtained from TargetScan to identify missense, non-missense, and amplification driver genes. After eliminating batch effects arising from shRNA thermodynamics and pooling, all denoised logFCs in each cell line were ranked. Then, 50% of the most consistent shRNAs of each gene were selected based on their Pearson correlation with the average shRNA profile. Finally, a negative binomial mixed effects model was applied to the shRNA ranks in three alteration classes: missense, non-missense, and amplification alterations. The random effect in this model is associated with the cell line and the fixed effect is a binary variable indicating mutational status.

Step 2: batch-effect removal and ranking

We then used the `removeBatchEffect` function in `limma` package in R³³ to remove two batch effects, namely pool and thermodynamic stability of seed sequences, from the shRNA data. Subsequently, we ranked all the corrected shRNA values per cell line such that shRNAs with lower ranks represented higher depletions.

Step 3: obtaining consistent shRNAs

Similar to the ATARiS algorithm, though with a different approach, we found a subset of shRNAs for each gene which exhibited consistent behavior across all cell lines. To achieve this, we calculated Pearson correlation coefficients between all shRNA viability ranks of a gene and the average viability ranks across all shRNAs for all cell lines in a given gene and retained half of the shRNAs with the highest Pearson correlation coefficients. The final preprocessed data to perform NBDep algorithm contains information on 76,495 shRNAs targeting 7241 genes 339 cell lines.

Step 4: statistical testing using a negative binomial mixed effects model

Finally, a negative binomial mixed effects model was performed to discover gene drivers in the refined data obtained in the previous steps. To incorporate the heterogeneity of cell lines into the model, they were entered as a random effect covariate in the model. In this model, we assumed that the viability rank for i th shRNA and k th cell line, denoted by R_{ik} , follows a negative binomial distribution

$$R_{ik} \sim NB(\mu_{ik}, \theta),$$

for $i = 1, 2, \dots, n$; $k = 1, 2, \dots, l$. Parameters μ_{ik} and θ are mean and shape parameter of negative binomial distribution. We selected the negative binomial distribution because it is a two-parameter distribution that offers the necessary flexibility to model non-negative integer values, making it suitable for the ranks. The probability mass function of R_{ik} is given by

$$f(r) = \frac{\Gamma(r + \theta)}{\Gamma(\theta)r!} \left(\frac{\theta}{\mu_{ik} + \theta} \right)^\theta \left(\frac{\mu_{ik}}{\mu_{ik} + \theta} \right)^r,$$

where $\Gamma(\cdot)$ is the gamma function. The parameter μ_{ik} is defined as:

$$\mu_{ik} = E(R_{ik}|X_{ik}; u_k) = \exp(X_{ik}\beta + u_k).$$

The mutational status of i th shRNA in k th cell line is a fixed effect, denoted as the variable X_{ik} where 0 and 1 indicate wild-type and mutant cell lines, respectively. The variable u_k represents the random effect for the k th cell line. Under the assumption of this model, the vector $u = (u_1, u_2, \dots, u_l)^T$ follows a multivariate normal with mean 0 and covariance matrix Σ (i.e., $u \sim MVN(0, \Sigma)$). The maximum likelihood estimates of the model parameters, (θ, β, Σ) , were obtained using the `glmmTMB` package in R. To test for the dependency between mutational status and viability score of a gene, we performed one-sided hypothesis testing with $H_0 : \beta = 0$ and $H_1 : \beta < 0$. Under the null hypothesis, β follows a normal distribution and consequently the p-value is straightforward to compute.

Simulation study

In order to evaluate the performance of NBDep and gene-level methods, we conducted two simulation studies to assess type I error and statistical power.

Simulation study: type I error analysis

In this simulation study, we selected a set of genes for type I error analysis, including known cancer genes *KRAS*, *NRAS*, *BRAF*, *TP53*, *CDK4*, *AXIN1*, *DHX9*, *TOP2A*, *RRAS2*, *COL1A1* as well as a few genes that are not known to be involved in cancer namely *UGT8*, *TRPV3*, *FUS*, *BRX1*, *NAN60*, *RPS6KA4*, *PTBP1*, *SPRR1A*, *MYLIP*, *LRCH4*, *DUSP26*, *CSF2RB*, *CBX4*, *DLD*, *MCRS1*, *PYGB*, *ATP11B*, and *BTNL2*. We assumed that there were m mutated and w wild-type cell lines with respect to a given gene. To generate 5000 random datasets, we performed the following steps:

- (1) We randomly selected n_m cell lines from the set of all cell lines and assigned them to the mutant group. We generated n_m from a uniform distribution between 2 and m .
- (2) We randomly selected n_w cell lines from the remaining cell lines and assigned them to the wild-type group. We generated n_w from a uniform distribution between 2 and w . This ensured that both groups had the same distribution, and the generated data was suitable for assessment of type I error.
- (3) We then tested for a relationship between mutation status and viability using various methods and calculated type I error as the fraction of rejected tests over the total 5000 tests for each method.

Simulation study: statistical power analysis

In this section, we focused on known cancer genes having more than five mutant cell lines, identified through all methods including APSiC on ATARiS and DEMETER2 scores, as well as NBDep. These genes, namely *KRAS*, *NRAS*, *PIK3CA*, *TP53*, *CTNNB1*, and *BRAF*, were used in analyzing statistical power. To this end, we first simulated datasets in which the distributions of two groups were different. We then subsampled n_m mutated cell lines and n_w wild-type cell lines, where n_m was between 2 and 10, and $n_w = qn_m$ for each gene. The parameter q represents the ratio of wild-type to mutant cell lines, and we considered values of $q = 1, 2, 3, 4$. For each gene,

n_m and q , we generated 800 datasets. The statistical power was then calculated as the ratio of rejected tests to the total number of tests for each setting.

Statistical testing on gene-level data

We applied the Wilcoxon rank sum test and APSiC, a ranked-based statistical approach based on Irwin-Hall and Bates distributions, to identify self-dependencies from perturbation screens²¹. These tests were conducted on gene-level scores obtained from ATARiS or DEMETER2 dependency scores, which are commonly used for aggregating shRNA viability scores targeting a gene into a single gene-level score. Since APSiC employs a one-sided hypothesis testing, we used a one-sided hypothesis testing in NBDep too to facilitate direct comparison with APSiC results on gene-level ATARiS and DEMETER2 scores. We discovered that the statistical power of the Wilcoxon rank sum test was significantly lower than that of APSiC, and as a result, we only included the results of APSiC in our analyses.

PPI enrichment

We conducted an enrichment analysis on the set of genes identified through our pan-cancer analysis using data from the STRING database, which contains information on protein–protein interactions. The goal of this analysis was to determine if the identified genes had a greater number of interactions with each other compared to a set of randomly selected genes³⁴.

Multiple testing correction

For multiple testing correction, we controlled the false discovery rate using Benjamini and Hochberg approach³⁵. The significance level was set to 0.05.

Results

Simulated studies

We first assessed the type I error of our method on synthetic data (see “Methods”). At the significance level of 5%, the estimated type I error of NBDep was 0.067, slightly higher than the expected 5% error. Next, we compared the statistical power of our method with those of APSiC statistical test on ATARiS and DEMETER2 gene-level dependency scores on three well-known missense driver genes (*KRAS*, *NRAS*, and *PIK3CA*) using subsampling approaches. To determine whether our method is more robust to low sample size than other methods, we generated 800 datasets for each sample size resulting from various numbers of mutant cell lines and values of q (see “Materials and methods” section). In general, the statistical power of the NBDep algorithm is higher than other methods for small sample sizes with $q = 1$ (Fig. 2) and other values of q in the analyzed genes (Supplementary Figs. S2, S3).

Influence of denoising and ranking on NBDep’s performance

In this section, we investigated the reason why the NBDep algorithm has enhanced statistical power. In particular, we considered whether it is due to denoising or ranking. To this end, we compared the following approaches: (1) a linear mixed-effects model on raw logFC of shRNA readouts, (2) a linear mixed-effects model on denoised logFC values, (3) a negative binomial mixed-effects model on rank data obtained from raw logFC values, (4) a negative binomial mixed-effects model on rank data obtained from denoised logFC values (i.e., NBDep algorithm). We performed a similar simulation as the previous section on *KRAS*, *NRAS*, *PIK3CA*. The Supplementary Fig. S4 demonstrates that the main reason for the enhanced statistical power is using ranks of logFC values. However, denoising can also slightly improve the statistical power. In order to ensure that the ranking does not result in elevated false positives, we calculated the type one error for the third model, which was 0.058 and plausible.

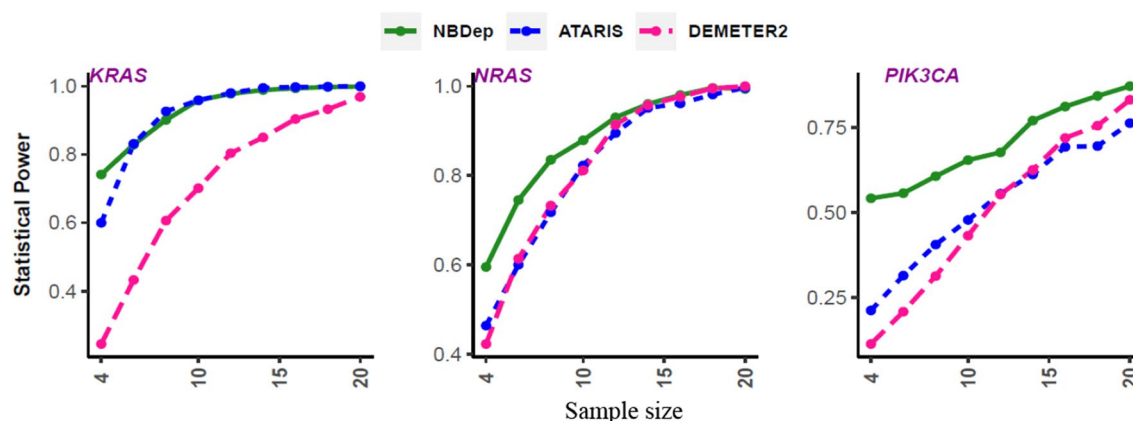


Figure 2. Comparison of statistical power between APSiC method applied to ATARiS and DEMETER2 scores and the NBDep algorithm for *KRAS*, *NRAS*, and *PIK3CA* across various sample sizes and $q = 1$.

Pan-cancer analysis

NBDep method identified 89 missense mutational driver genes, 41 amplification driver genes, and 125 non-missense mutational driver genes in pan-cancer analyses. Top identified amplification, missense, and non-missense mutational genes by the NBDep method in pan-cancer are depicted in Fig. 3c.

In analysis of missense driver genes, NBDep identified sixteen well-recognized genes reported in IntOGen namely *KRAS*, *NRAS*, *BRAF*, *PIK3CA*, *CTNNB1*, *TP53*, *SMAD4*, *BCL2*, *SIX1*, *FBXO11*, *COL1A1*, *MAP2K1*, *TOP2A*, *DHX9*, *HRAS* (Fig. 3a). In addition, unreported genes in IntOGen such as *CANT1*, *IL20RA*, *BRX1*, *BUB1*, *ATIC*, *Twist2*, and *CSF2RB* were identified by the NBDep method as putative missense mutational cancer genes

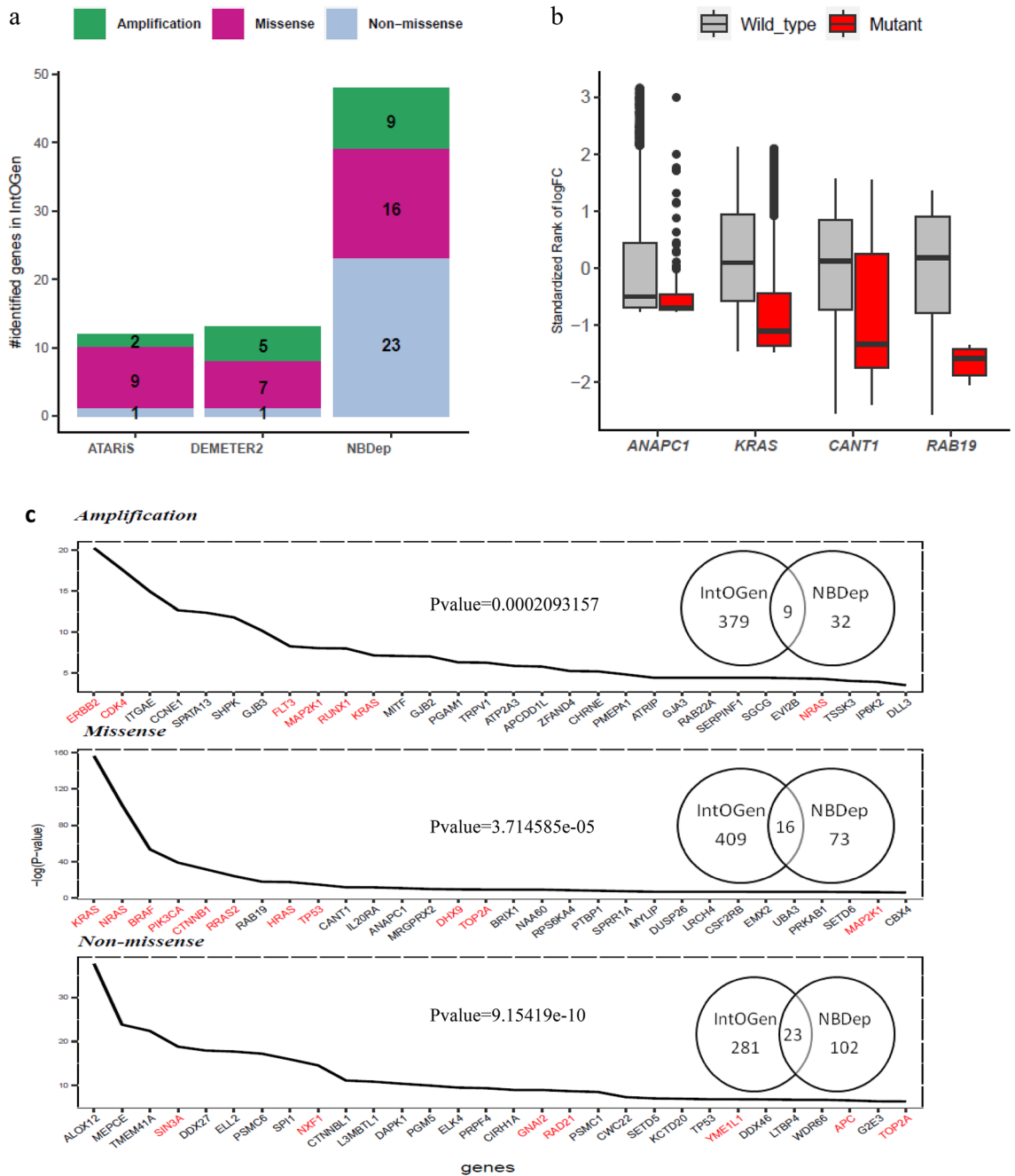


Figure 3. (a) The number of identified genes common with IntOGen in pan-cancer analysis resulted from three methods: NBDep as well as applying the APSiC method on ATARiS and DEMETER2 gene-level scores in three alterations: amplification, missense, and non-missense. (b) Box plots of standardized denoised ranks of one novel missense-driver gene, *ANAPC1*, one experimentally curated missense driver gene, *KRAS*, two genes previously reported in literatures, *RAB19* and *IL20RA*, identified by NBDep. (c) Top identified amplification, missense, and non-missense mutational genes by NBDep method in pan-cancer. Red color indicates IntOGen genes. Venn diagrams indicates the number of common genes with IntOGen in each of three alteration groups along with their enrichment p-values calculated by the hypergeometric test.

but these genes were not detected by using APSiC on ATARiS and DEMETER2 gene-level scores. Recently, the important role of these genes as driver genes in different stages of cancer has been presented^{36–42}. Box plots depicting the standardized, denoised ranks of four example genes identified by the NBDep algorithm are displayed in Fig. 3b. As depicted in Supplementary Fig. S5, there was a significant PPI enrichment observed for both the identified missense driver genes and the identified genes not reported in IntOGen, with p-values of 0.000214 and 0.00767, respectively. One of the top missense driver genes identified by NBDep method is *ANAPC1* which has not been reported in IntOGen and the literature research and it can be introduced as a novel cancer gene. Figure 4b indicates that *ANAPC1* has significantly lower ranks in missense cell lines than wild-type ones even before taking account for cell line effect in the negative binomial mixed effects model. On the other hand, as depicted in Fig. 4a, nine identified amplification driver genes by NBDep are in the IntOGen gene list, namely *ERBB2*, *CDK4*, *FLT3*, *MAP2K1*, *KRAS*, *NRAS*, *GNAS*, *RUNX1* as oncogenes and *LATS2* as a tumor suppressor. Additionally, other amplification driver genes identified by NBDep, such as *ITGAE*, *ALKBH3*, *SHPK*, and *GJB3* have been previously reported to play a significant role in cancer^{43–46}. NBDep also identified 23 non-missense driver genes that were previously reported in IntOGen, namely *SIN3A*, *NTRK1*, *FN1*, *MGA*, *MAP3K1*, *PTEN*, *DROSHA*, *WRN*, *RELA*, *CDX2*, *KDR*, *CDKN1B*, *POLD1*, *PIK3R1*, *NUP214*, *PML*, *GATA3*, *TOP2A*, *APC*, *TP53*,

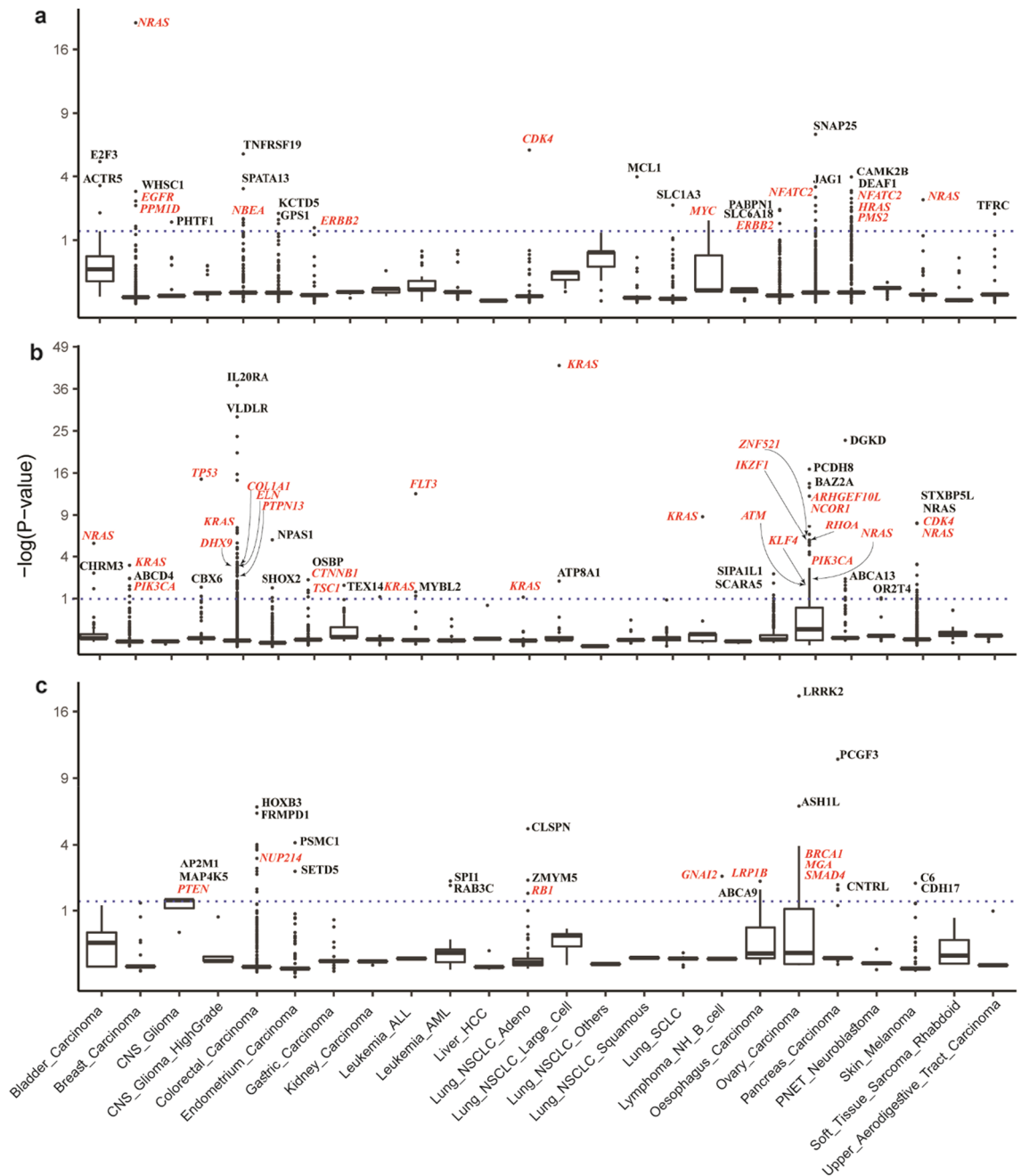


Figure 4. The identified amplification, missense, and non-missense driver genes by NBDep in cancer-specific analyses. The genes highlighted in red are in the IntOGen list.

RAD21, *GNAI2*, and *NXF1*. Recent studies have also highlighted the role of some of the top non-missense driver genes that were not previously identified by IntOGen, such as *ALOX12*⁴⁷. Furthermore, NBDep identified two putative non-missense driver genes *MEPCE* and *TMEM41A* that were not previously known to be cancer driver genes.

In the pan-cancer analysis of all three alterations, several genes were identified as both non-missense and missense driver genes, including *AT1C*, *MCRS1*, *TBL3*, *CACNA1A*, *TOP2A*, and *TP53*. Moreover, *NRAS*, *KRAS*, and *MAP2K1* were found to be amplification and missense driver genes. Additionally, *KCNQ1* was identified as a driver gene in cancer by NBDep but has not been previously reported in both non-missense and amplification driver genes. Figure 4a shows that the NBDep algorithm identified more common genes across all three alterations compared to using the APSiC method on ATARiS and DEMETER2 scores. Figure 4c displays the top 30 genes identified by the NBDep algorithm in all three alterations, along with their p-values computed using the hypergeometric test for finding the statistical significance of overlaps between identified genes and the IntOGen gene list. It is worth noting that the number of identified driver genes by NBDep was higher in cancer types with more cell lines such as colorectal, skin, and pancreatic carcinoma.

Cancer-specific analysis

We employed the NBDep algorithm to identify driver genes across 26 cancer types only on genes with at least two mutant and two wild-type cell lines. Our analysis revealed 87 amplification, 177 missense, and 45 non-missense driver genes across cancer types. The most frequent identified amplification driver genes belonged to pancreatic and ovarian carcinoma. The identified missense mutational genes were found in colorectal carcinoma and ovarian carcinoma, and non-missense mutational genes were observed in colorectal carcinoma.

Among the amplification driver genes, *NRAS* (in breast carcinoma and skin melanoma), *ERBB2* (in gastric carcinoma, esophageal carcinoma), and *NFATC2* (in ovarian carcinoma and pancreatic carcinoma) were identified in more than one cancer type. Four genes, namely *KRAS* (in breast carcinoma, non-small cell lung cancer (NSCLC), lymphoma multiple myeloma carcinoma, colorectal carcinoma, and leukemia), *NRAS* (in pancreatic carcinoma, bladder carcinoma, primitive neuro-ectodermal tumors (PNET)), *PIK3CA* (in breast carcinoma and ovarian carcinoma), and *ARHGAP31* (in colorectal carcinoma and pancreatic carcinoma) were identified as missense driver genes in more than one cancer type. Notably, the number of identified driver genes by NBDep was higher in cancer types with more cell lines, such as colorectal, skin, and pancreatic carcinoma.

Among the 89 identified missense driver genes in pan-cancer, 21 genes were detected in multiple cancer types, including known cancer genes such as *TP53*, *KRAS*, *NRAS*, and *PIK3CA*. Although most of these genes were identified in colorectal cancer, *PIK3CA*, *TP53* and *CTNNB1* were detected in other cancer types too (breast carcinoma, central nervous system glioma high grade, and gastric carcinoma). *KRAS*, *IL2ORA*, *CANTI*, *NRAS* and *TP53* were the top genes in different cancers. A recent study indicates *IL2ORA* is an important regulator of oncogenic and immune pathways in colorectal carcinoma³⁶ (Supplementary Table 1).

NRAS (breast carcinoma and skin melanoma), *CDK4* (NSCLC), and *ERBB2* (gastric carcinoma and esophageal carcinoma) were identified as the top identified amplification driver genes in different cancer types (Supplementary Table 1).

Also, we found three top genes which have not been reported in IntOGen, *TNFRSF19*, *GJA3*, and *GJB2*, recently found by researchers to have main role in colorectal cancer^{48,49} (Supplementary Table 1).

Our method introduced *SPATA13* a guanine-factor as a novel amplification driver gene in colorectal carcinoma, required for *MMP9* up-regulation via the JNK signaling pathway in colorectal tumor cells. Also, NBDep method identified amplification *WHSC1* in breast carcinoma as a novel gene.

Additionally, 16 identified non-missense mutational driver genes in the pan-cancer analyses were also detected as non-missense genes across cancer types, including *CLSPN* and *PIK3C3* (NSCLC), *PSMC1* and *SETD5* (endometrial carcinoma), *MGA* (ovarian carcinoma), *PTEN* (CNS glioma), *SP11* (leukemia), *GNAI2* (lymphoma NH B-cell), and the eight genes in colorectal carcinoma.

For the cancer types with small number of cell lines, the NBDep method is able to identify well-known genes. For leukemia, lymphoma NH B-cell, lymphoma multiple myeloma cancer, and NSCLC that have very low number of cell lines, *KRAS* was detected as the most important missense gene driver. NBDep found *MYC* as an amplification driver gene in lymphoma multiple myeloma cancer while gene-level approaches did not identify this gene. Our finding was confirmed by a recent study where amplified *MYC* was shown to be effective in myeloma cancer⁵⁰. In addition, the recent research approves the identified *ATP8A1* by NBDep as an important gene in NSCLC⁵¹. The NBDep method also identified 43 missense driver genes in ovarian cancer having a significant protein–protein network (p-value = 2.12e–06) where well-known genes such as *RHOA*, *PIK3CA*, and *ATM* are among these genes^{52–54}. Our method proposed *PCDH8* and *BAZ2A* as novel missense driver genes in ovarian carcinoma.

In summary, the NBDep algorithm is able to identify well-known driver genes in cancer-specific analyses with large and small number of cell lines and to introduce novel putative driver genes (Fig. 4). In addition, NBDep algorithm identified more curated genes regarding to IntOGen in cancer-specific than ATARiS and DEMETER2 (Supplementary Fig. S6).

Discussion

Perturbation screens including RNAi screening has become increasingly popular in the field of cancer genomics over the past decade. One major limitation of RNAi screening is the off-targets issue where it poses a major challenge to infer actual gene effects in these screens. Different computational methods, including ATARiS and DEMETER2, have been developed to handle off-target effects of shRNAs, leading to gene-level scores known as dependency scores. ATARiS achieves an aggregate score for each gene by discarding shRNAs with non-consistent

behavior across all cell lines. DEMETER2 uses a hierarchical Bayesian model to explicitly handle off-target effect associated with seed sequence of each shRNA resulting to gene-level dependency scores.

In this research, we presented a statistical framework aimed at handling off-target effects at shRNA level to identify driver genes associated to missense, amplification, and non-missense alterations. We applied our approach to 26 cancer-specific types as well as pan-cancer data using the Project DRIVE data. We coped with off target effects of shRNAs by incorporating thermodynamic stability of 7-mer seed shRNAs at the shRNA level as a main batch effect proposed by TargetScan. After calculating the denoised shRNA from the original shRNA scores, we then employed 50% of consistent shRNAs of each gene across cell lines. We subsequently performed a negative binomial mixed effects model to investigate association of gene perturbation and alteration statuses of genes. We additionally showed that NBDep algorithm is robust in small sample sizes and can detect driver genes more effectively than using the APSiC method on ATARiS and DEMETER2 gene-level scores.

Having compared the numbers of shRNAs designed for genes in IntOGen to other genes, we recognized that IntOGen genes have more shRNAs than non-well-known genes. This finding advocates that it may be more efficient to design more reagents for all genes when performing RNAi screening.

Our method was capable of identifying well-known genes in pan-cancer and 26 cancer-specific types including *KRAS*, *NRAS*, *TP53*, *BRAF*, and *PIK3CA* in missense alteration, *ERBB2*, *CDK4*, and *FLT3* as amplification driver genes, and well-recognized non-missense driver genes such as *TP53*, *TOP2A*, *APC*, and *GATA3*. Additionally, in cancer-specific analyses, our method was able to identify well-known genes such as *KRAS* and *PIK3CA* as missense driver genes, *NRAS* and *EGFR* as amplification driver genes in breast carcinoma. NBDep method could not identify any driver gene in liver carcinoma. NBDep identified *CCKAR* and *WHSC* as novel missense and amplification driver genes, respectively, in breast carcinoma. *CCKAR* was reported as a driver gene in gall-bladder and biliary tract cancer and *WHSC* as *TP53* binding protein⁵⁵. Moreover, *PCDH8*, a protein coding that acts as a cell adhesion molecule, and *BAZ2A* having DNA binding activity, were suggested as missense novel genes in ovarian cancer.

In summary, our method utilizes more information of shRNAs in RNAi screening and is capable to identify above-mentioned self-dependencies in pan-cancer and 26 cancer-specific types while handling off-target effects. Our approach also provides denoised shRNA ranks through which it is possible to explore other types of cancer dependencies such as synthetic lethality.

Data availability

The raw shRNA data and ATARiS scores for Project DRIVE were obtained from the Mendeley Data repository (version 4) at <https://data.mendeley.com/datasets/y3ds55n88r>, version 4. The DEMETER2 dependency scores from Project DRIVE were retrieved from the DepMap project (version 22Q4). Additionally, molecular profiling data was acquired from the DepMap project (version 21Q4). The CCLE GISTIC copy number alteration data was obtained from the CBioPortal website (Cancer Cell Line Encyclopedia, Novartis/Broad, Nature 2012). The thermodynamic stability of the seed was extracted from supplementary data 5 of the original paper.

Received: 2 August 2023; Accepted: 5 January 2024

Published online: 09 January 2024

References

- Sundara Rajan, S., Ludwig, K. R., Hall, K. L., Jones, T. L. & Caplen, N. J. Cancer biology functional genomics from small RNAs to big dreams. *Mol. Carcinog.* **59**, 1343–1361 (2020).
- Knott, G. J. & Jennifer, A. CRISPER-Cas guides the future of genetic engineering. *Science* **1361**, 866–869 (2018).
- Dai, M. *et al.* In vivo genome-wide CRISPR screen reveals breast cancer vulnerabilities and synergistic mTOR/hippo targeted combination therapy. *Nat. Commun.* **12**, 3055 (2021).
- Arfaoui, A. *et al.* A genome-wide RNAi screen reveals essential therapeutic targets of breast cancer stem cells. *EMBO Mol. Med.* **11**, e9930 (2019).
- Abdelrahim, M., Safe, S., Baker, C. & Abudayyeh, A. RNAi and cancer implications and applications. *J. RNAi Gene Silenc.* **2**(1), 136 (2006).
- Chun-Qing, S. *et al.* Genome-wide CRISPR screen identifies regulators of mitogen-activated protein kinase as suppressors of liver tumors in mice. *Gastroenterology* **152**(5), 1161–1173 (2017).
- Amanda, B., Andreas, K. & Karol, K. *RNAi and Off-Target Effects* 3–20 (Bentham Science Publisher, 2014).
- Brown, K. & Samarsky, D. RNAi off-targeting: Light at the end of the tunnel. *J. RNA Gene Silenc.* **2**, 175–177 (2006).
- Shao, D. D. *et al.* ATARiS: Computational quantification of gene suppression phenotypes from multisample RNAi screens. *Genome Res.* **23**, 665–678 (2013).
- Schmich, F. *et al.* gesper: A statistical model for deconvoluting off-target-confounded RNA interference screens. *Genome Biol.* **16**, 220 (2015).
- McFarland, J. M. *et al.* Improved estimation of cancer dependencies from large-scale RNAi screens using model-based normalization and data integration. *Nat. Commun.* **9**, 4610 (2018).
- Konig, R. *et al.* A probability-based approach for the analysis of large-scale RNAi screens. *Nat. Methods* **4**, 847–849 (2007).
- Tsherniak, A. *et al.* Defining a cancer dependency map. *Cell* **170**, 564–576.e16 (2017).
- Agarwal, V., Bell, G. W., Nam, J. W. & Bartel, D. P. Predicting effective microRNA target sites in mammalian mRNAs. *Elife* **4**, e05005 (2015).
- Shimomura, I., Yamamoto, Y. & Ochiya, T. Synthetic lethality in lung cancer—from the perspective of cancer genomics. *Medicines* **6**(1), 38 (2019).
- Dolly, S. O. *et al.* RNAi screen reveals synthetic lethality between cyclin G-associated kinase and FBXW7 by inducing aberrant mitoses. *Br. J. Cancer* **117**, 954–964 (2017).
- Maia, A. F. *et al.* Genome-wide RNAi screen for synthetic lethal interactions with the *C. elegans* kinesin-5 homolog BMK-1. *Sci. Data* **2**, 150020 (2015).
- Srivatsa, S. *et al.* Discovery of synthetic lethal interactions from large-scale pan-cancer perturbation screens. *Nat. Commun.* **13**, 7748 (2022).
- O’Neil, N., Bailey, M. & Hieter, P. Synthetic lethality and cancer. *Nat. Rev. Genet.* **18**(10), 613–623 (2017).

20. McDonald, E. R. *et al.* Project DRIVE: A compendium of cancer dependencies and synthetic lethal relationships uncovered by large-scale, deep RNAi screening. *Cell* **170**, 577–592 (2017).
21. Montazeri, H. *et al.* Systematic identification of novel cancer genes through analysis of deep shRNA perturbation screens. *Nucleic Acids Res.* **49**, 8488–8504 (2021).
22. Rameseder, J. *et al.* A multivariate computational method to analyze high-content RNAi screening data. *J. Biomol. Screen.* **20**, 985–997 (2015).
23. Luo, B. *et al.* Highly parallel identification of essential genes in cancer cells. *Proc. Natl. Acad. Sci.* **105**, 20380–20385 (2008).
24. Li, W. *et al.* MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9 knockout screens. *Genome Biol.* **15**, 1–12 (2014).
25. Suomi, T., Corthals, G. L., Nevalainen, O. S. & Elo, L. L. Using peptide-level proteomics data for detecting differentially expressed proteins. *J. Proteom. Res.* **14**, 4564–4570 (2015).
26. Goeminne, L. J., Argentini, A., Martens, L. & Clement, L. Summarization vs peptide-based models in label-free quantitative proteomics: Performance, pitfalls, and data analysis guidelines. *J. Proteom. Res.* **14**, 2457–2465 (2015).
27. Goeminne, L. E., Gevaert, K. & Clement, L. Peptide-level robust ridge regression improves estimation, sensitivity, and specificity in data-dependent quantitative label-free shotgun proteomics. *Mol. Cell. Proteom.* **15**, 657–668 (2016).
28. Marcotte, R. *et al.* Functional genomic landscape of human breast cancer drivers, vulnerabilities, and resistance. *Cell* **164**, 293–309 (2016).
29. Yu, J., Silva, J. & Califano, A. ScreenBEAM: A novel meta-analysis algorithm for functional genomics screens via Bayesian hierarchical modeling. *Bioinformatics* **32**, 260–267 (2016).
30. Mermel, C. H. *et al.* GISTIC2.0 facilitates sensitive and confident localization of the targets of focal somatic copy-number alteration in human cancers. *Genome Biol.* <https://doi.org/10.1186/gb-2011-12-4-r41> (2011).
31. Garcia, D. M. *et al.* Weak seed-pairing stability and high target-site abundance decrease the proficiency of lsy-6 and other micro-RNAs. *Nat. Struct. Mol. Biol.* **18**, 1139–1146 (2011).
32. Martínez-Jiménez, F. *et al.* A compendium of mutational cancer driver genes. *Nat. Rev. Cancer* **20**, 555–572 (2020).
33. Ritchie, M. E. *et al.* Limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
34. Szklarczyk, D. *et al.* The STRING database in 2021: Customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Res.* **49**, D605–D612 (2021).
35. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B* **57**, 289–300 (1995).
36. Yu, D. *et al.* Super-enhancer induced IL-20RA promotes proliferation/metastasis and immune evasion in colorectal cancer. *Front. Oncol.* **11**, 724655 (2021).
37. Ge, J., Huang, X., Wang, P. & Lu, C. Expression of biogenesis of ribosomes BRX1 is associated with malignant progression and prognosis in colorectal cancer. *Transl. Cancer Res.* **9**, 5595–5602 (2020).
38. Gerhardt, J. *et al.* The androgen-regulated calcium-activated nucleotidase 1 (CANT1) is commonly overexpressed in prostate cancer and is tumor-biologically relevant in vitro. *Am. J. Pathol.* **178**, 1847–1860 (2011).
39. Jiang, N. *et al.* BUB1 drives the occurrence and development of bladder cancer by mediating the STAT3 signaling pathway. *J. Exp. Clin. Cancer Res.* **40**, 1–17 (2021).
40. Niu, N. *et al.* ATIC facilitates cell growth and migration by upregulating Myc expression in lung adenocarcinoma. *Oncol. Lett.* **23**(4), 1–11 (2022).
41. Fang, X. *et al.* Twist2 contributes to breast cancer progression by promoting an epithelial-mesenchymal transition and cancer stem-like cell self-renewal. *Oncogene* **30**, 4707–4720 (2011).
42. Rashid, M. *et al.* Discovery of a novel potentially transforming somatic mutation in CSF2RB gene in breast cancer. *Cancer Med.* **10**, 8138–8150 (2021).
43. Tasaki, M., Shimada, K., Kimura, H., Tsujikawa, K. & Konishi, N. ALKBH3, a human AlkB homologue, contributes to cell survival in human non-small-cell lung cancer. *Br. J. Cancer* **104**(4), 700–706 (2011).
44. Franceschi, S. *et al.* Sedoheptulose kinase SHPK expression in glioblastoma: Emerging role of the nonoxidative pentose phosphate pathway in tumor proliferation. *Int. J. Mol. Sci.* **23**, 5978 (2022).
45. Huo, Y. *et al.* GJB3 promotes pancreatic cancer liver metastasis by enhancing the polarization and survival of neutrophil. *Front. Immunol.* **13**, 983116 (2022).
46. Hu, X. *et al.* ITGAE defines CD8+ tumor-infiltrating lymphocytes predicting a better prognostic survival in colorectal cancer. *EBioMedicine* **35**, 178–188 (2018).
47. Huang, Z., Xia, L., Zhou, X., Wei, C. & Mo, Q. ALOX12 inhibition sensitizes breast cancer to chemotherapy via AMPK activation and inhibition of lipid synthesis. *Biochem. Biophys. Res. Commun.* **514**(1), 24–30 (2019).
48. Liu, Y. J. *et al.* An analysis regarding the association between connexins and colorectal cancer (CRC) tumor microenvironment. *J. Inflamm. Res.* **15**, 2461–2476 (2022).
49. Schön, S. *et al.* β -catenin regulates NF- κ B activity via TNFRSF19 in colorectal cancer cells. *Int. J. Cancer* **135**, 1800–1811 (2014).
50. Holien, T. *et al.* MYC amplifications in myeloma cell lines: Correlation with MYC-inhibitor efficacy. *Oncotarget* **6**, 22698–22705 (2015).
51. Li, D. *et al.* The role of ATP8A1 in non-small cell lung cancer. *Int. J. Clin. Exp. Pathol.* **10**, 7760–7766 (2017).
52. Wang, X., Jiang, W., Kang, J., Liu, Q. & Nie, M. Knockdown of RhoA expression alters ovarian cancer biological behavior in vitro and in nude mice. *Oncol. Rep.* **34**, 891–899 (2015).
53. Thorstenson, Y. R. *et al.* Contributions of ATM mutations to familial breast and ovarian cancer. *Cancer Res.* **63**, 3325–3333 (2003).
54. Campbell, I. G. *et al.* Mutation of the PIK3CA gene in ovarian and breast cancer. *Cancer Res.* **64**, 7678–7681 (2004).
55. Xu, H. L. *et al.* Variants in CCK and CCKAR genes to susceptibility to biliary tract cancers and stones: A population-based study in Shanghai, China. *J. Gastroenterol. Hepatol.* **28**, 1476–1481 (2013).

Author contributions

Z.T. and H.M. designed the research project. Z.T. developed the proposed method and conducted the analyses, incorporating feedback from H.M. Z.T. and H.M. wrote the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-024-51453-5>.

Correspondence and requests for materials should be addressed to H.M.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024