

Network perturbation analysis of gene transcriptional profiles reveals protein targets and mechanism of action of drugs and influenza A viral infection

Heeju Noh^{1,2}, Jason E. Shoemaker^{3,4} and Rudiyanto Gunawan^{1,2,*}

¹Institute for Chemical and Bioengineering, ETH Zurich, Zurich 8093, Switzerland, ²Swiss Institute of Bioinformatics, Lausanne 1015, Switzerland, ³Department of Chemical and Petroleum Engineering, University of Pittsburgh, Pittsburgh, PA 15261, USA and ⁴Department of Computational and Systems Biology, University of Pittsburgh, Pittsburgh, PA 15261, USA

Received October 13, 2017; Revised December 15, 2017; Editorial Decision December 22, 2017; Accepted December 22, 2017

ABSTRACT

Genome-wide transcriptional profiling provides a global view of cellular state and how this state changes under different treatments (e.g. drugs) or conditions (e.g. healthy and diseased). Here, we present ProTINA (Protein Target Inference by Network Analysis), a network perturbation analysis method for inferring protein targets of compounds from gene transcriptional profiles. ProTINA uses a dynamic model of the cell-type specific protein–gene transcriptional regulation to infer network perturbations from steady state and time-series differential gene expression profiles. A candidate protein target is scored based on the gene network’s dysregulation, including enhancement and attenuation of transcriptional regulatory activity of the protein on its downstream genes, caused by drug treatments. For benchmark datasets from three drug treatment studies, ProTINA was able to provide highly accurate protein target predictions and to reveal the mechanism of action of compounds with high sensitivity and specificity. Further, an application of ProTINA to gene expression profiles of influenza A viral infection led to new insights of the early events in the infection.

INTRODUCTION

The identification of the molecular targets of pharmacologically relevant compounds is vital for understanding the mechanism of action (MoA) of drugs, as well as for exploring off-target effects. While the definition of a target can be quite arbitrary, the term generally refers to a molecule whose interaction with the compound is connected to the compound’s effects (1). In this study, transcription factors (TFs) and their protein interaction partners represent the

target molecules, while differential gene expression profiles represent the effects. Among existing technologies for protein target discovery (e.g. biochemical affinity purification, RNAi knockdown or gene knockout experiments) (2), gene expression profiling has received much recent attention due to its relative ease of implementation as well as the availability of large-scale public databases and well-established experimental protocols and data analytical methods. A complication when using gene expression profiling for target discovery is that the data give only indirect indications of the drug’s action. As illustrated in Figure 1A, the interaction between a compound and its protein target(s) is expected to result in the differential expression of downstream genes that are regulated by the protein target(s). But, the expression of the protein targets themselves may not—and often do not—change (3). Consequently, target discovery using gene expression profiles requires computational methods to identify the (upstream) targets from the (downstream) effects.

Existing computational strategies for compound target identification using gene expression profiles can generally be classified into two groups: comparative analysis and network-based analysis (4). Comparative analysis methods use the gene expression profiles as drug signatures. Here, the similarity between the differential gene expression of a drug treatment and those of reference compounds or experiments with known targets, implies a closeness in the molecular targets and the MoA. A notable example of such an approach is the Connectivity Map (5), which provides gene expression profiles of human cell lines treated by ~5000 small molecule compounds as queryable signatures for evaluating drug–drug similarities (6). The obvious drawback of comparative analysis methods is their dependence on an extensive and accurate target annotation of the reference gene expression profiles.

In network-based analysis, one adopts a system-oriented view by using cellular networks, such as gene regulatory network (GRN) and/or protein–protein interaction network

*To whom correspondence should be addressed. Tel: +41 44 633 21 34; Fax: +41 44 633 12 52; Email: rudi.gunawan@chem.ethz.ch

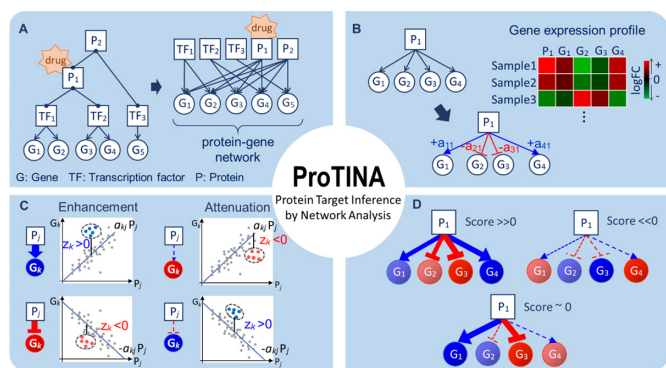


Figure 1. Protein target prediction by ProTINA. (A) The protein–gene network describes direct and indirect regulations of gene expression by transcription factors (TF) and their protein partners (P), respectively. A drug interaction with a protein is expected to cause differential expression of the downstream genes in the PGN. (B) Based on a kinetic model of gene transcriptional process, ProTINA infers the weights of the protein–gene regulatory edges, denoted by a_{kj} , using gene expression data. The variable a_{kj} describes the regulation of protein j on gene k , where the magnitude and sign of a_{kj} indicate the strength and mode ($+a_{kj}$: activation, $-a_{kj}$: repression) of the regulatory interaction, respectively. (C) A candidate protein target is scored based on the deviations in the expression of downstream genes from the PGN model prediction (P_j : $\log_2\text{FC}$ expression of protein j , G_k : $\log_2\text{FC}$ expression of gene k). The colored dots in the plots illustrate the $\log_2\text{FC}$ data of a particular drug treatment, while the lines show the predicted expression of gene k by the (linear) PGN model. The variable z_k denotes the z -score of the deviation of the expression of gene k from the PGN model prediction. A drug-induced enhancement of protein–gene regulatory interactions is indicated by a positive (negative) z_k in the expression of genes that are activated (repressed) by the protein (i.e. $a_{kj}z_k > 0$). Vice versa, a drug-induced attenuation is indicated by a negative (positive) z_k in the expression of genes that are activated (repressed) by the protein (i.e. $a_{kj}z_k < 0$). (D) The score of a candidate protein target is determined by combining the z -scores of the set of regulatory edges associated with the protein in the PGN. A positive (negative) score indicates a drug-induced enhancement (attenuation). The larger the magnitude of the score, the more consistent is the drug induced perturbations (enhancement/attenuation) on the protein–gene regulatory edges.

(PIN). A number of network-based analytical methods relied on dynamic models of the GRN to infer network perturbations caused by drug treatments (7–11). Several notable methods include Network Identification by multiple Regression (NIR) (7), Mode of action by Network Identification (MNI) (8), Sparse Simultaneous Equation Model (SSEM) (9) and DeltaNet (10,11). In these methods, the GRN is inferred from a training dataset of gene expression profiles using a linear regression derived from a dynamic mechanistic model of the gene transcriptional process. Subsequently, the inferred GRN is utilized for target identification to evaluate deviations in the differential gene expression caused by drug treatments (7–11) or in disease (12). A major pitfall of the above methods is that the inference of GRN from gene transcriptional profiles is highly challenging (13), as the inference problem often becomes underdetermined (i.e. the GRN may not be inferable) (14,15). In addition, as mentioned above, the expressions of the drug targets are often unaffected by the drug treatment (3).

Another group of network-based analytical methods utilizes cellular network graphs, either curated from the literature knowledge or inferred from gene expression data, to formulate statistical hypothesis tests for ranking drug

targets (3,16–22). Several methods in this category prioritize targets based on the enrichment of the downstream or neighbouring molecules in the network for differentially expressed genes, following the principle of ‘guilt by association’ (3,16–20). Another set of methods rank targets by scoring hypotheses that are generated based on causal relationships in the biological networks (21,22). A recent method called Detecting Mechanism of Action by Network Dysregulation (DeMAND), combines the GRN and PIN to create a molecular interaction network, where the drug targets are scored based on the statistical significance of drug-induced alterations in the joint gene expression distribution between two connected genes in the network (23). The methods in this group make use of only the (static) topology of cellular networks without much consideration of the dynamics of the gene transcriptional process, and thus are unable to fully exploit information contained in time-series datasets.

In this work, we developed ProTINA (Protein Target Inference by Network Analysis), a network perturbation analysis method for protein target identification using gene transcriptional profiles. The analysis involves two key steps: (a) the creation of a model of tissue or cell type-specific protein–gene regulatory network (PGRN) and (b) the calculation of protein target scores based on the enhancement or attenuation of the protein–gene regulations. In developing ProTINA, we addressed some of the drawbacks in the existing methods. First, the PGRN in ProTINA is based on a dynamic model of the gene transcriptional process, and is therefore able to take advantage of time-series gene expression profiles that are commonly generated by drug treatment studies. In addition, ProTINA leverages on the availability of comprehensive maps of protein–protein and protein–DNA interactions for the construction of the PGRN, which serves as prior information to alleviate network inferability issue. Finally, ProTINA scores candidate targets based on drug-induced perturbations to the expression of genes regulated by the targets, rather than the expression of the targets themselves. We demonstrated the superiority of ProTINA over the state-of-the-art method DeMAND and differential gene expression analysis (DE), in predicting the protein targets of drugs using human and mouse datasets from NCI-DREAM drug synergy challenge (24), genotoxicity study (25) and chromosome drug targeting study (26). Besides protein targets of compounds, we presented the application of ProTINA to study host-pathogen interactions, specifically for elucidating the targets of influenza A viral proteins.

MATERIALS AND METHODS

Gene expression data

We applied ProTINA to three datasets of drug treatments from NCI-DREAM drug synergy challenge (24), genotoxicity study (25) and chromosome drug targeting study (26), and to gene expression data of human lung cancer cell Calu-3 from influenza A viral infection studies (27–30). For NCI-DREAM drug synergy challenge, we obtained the raw *Affymetrix Human Genome U219* microarray data from Gene Expression Omnibus (GEO) database (31) (accession number: GSE51068). The raw data were first nor-

malized and transformed into \log_2 -scaled expressions using *justRMA* function in the *affy* package of Bioconductor (32). Then, the \log_2 fold change (\log_2 FC) differential expressions and their statistical significance (Benjamini–Hochberg adjusted *P*-values) were calculated using a linear fit model and empirical Bayes method in the *limma* package of Bioconductor. Three samples from the drug treatment using the low dose of Aclacinomycin A were dropped because all of the \log_2 FC expressions were close to 1 and thus not statistically significant. The probe sets were mapped to gene symbols using *hgu219.db* annotation package (Entrez Gene database as of 27 September 2015). In the case of multiple probe sets mapping to a gene symbol, we assigned the \log_2 FC from the probe set with the smallest average adjusted *P*-value over the samples.

The raw microarray data from genotoxicity study (25) in human HepG2 cell line were obtained from GEO (accession numbers: GSE28878 using *Affymetrix GeneChip Human Genome U133 Plus 2.0* array and GSE58235 using *Affymetrix HT Human Genome U133+ PM* array). As with the drug synergy data, the microarray data were first normalized using *justRMA*, and the \log_2 FCs and their adjusted *P*-values were calculated using *limma* in Bioconductor. Because the data came from different microarray platforms, the gene symbols were matched separately for each platform using *hgu133plus2.db* annotation package (Entrez database of 27 September 2015) and *HT_HG-U133_Plus_PM* annotation file in Affymetrix, respectively. Likewise, in the case of multiple probe sets matching a gene symbol, the probe set with the smallest average adjusted *P*-value across all samples was chosen.

The raw data from the chromosome-targeting study using mouse pancreatic alpha and beta cells (26) were also obtained from GEO database (accession number: GSE36379). Again, the raw data were normalized using *justRMA*, and the \log_2 FCs and their adjusted *P*-values were calculated by *limma*. The probes were mapped to the corresponding gene symbols using *moe430a.db* package (Entrez Gene database as of 27 September 2015) in Bioconductor. In the case of multiple probe sets mapping to a gene symbol, we selected the probe set with the smallest average adjusted *P*-value among the samples.

For influenza A infection analysis, we obtained the raw microarray data of four influenza studies (27–30) from GEO database (accession numbers: GSE40844, GSE37571, GSE33142 and GSE28166). The raw data were background-corrected and normalized using *normexp* and *quantile* methods in *limma* package of Bioconductor. The \log_2 FCs and their adjusted *P*-values were again calculated by *limma*. The probes were mapped to the corresponding gene symbols using *hgug4112a.db* package (Entrez Gene database as of 27 September 2015). Like before, for genes with multiple probe sets, we chose the \log_2 FC value corresponding to the probe set with the smallest average adjusted *P*-value.

Protein target identification using ProTINA

Protein–gene regulatory network. In ProTINA, the PGRN is a bipartite graph with weighted, directed edges pointing from a protein to a gene (see Figure 1A). The edges

in the PGRN describe the regulation of gene expression by TFs and their protein partners, the molecular targets of interest in this work. The bipartite PGRN above is able to capture feedback loops in the gene transcriptional regulation, even though these loops are not drawn explicitly. An example of such a feedback loop is when a protein regulates the expression of its own transcription factor(s). The PGRN is constructed by combining two types of networks, namely the TF–gene network and PIN. For the construction of human cell type-specific PGRNs, we relied on the Regulatory Circuit resource that provides 394 cell type and tissue-specific TF–gene interactions (33). More specifically, for the analysis of the NCI-DREAM drug synergy, genotoxic compound study, and influenza A viral infection study datasets, we used the TF–gene networks of human lymphoma cells, pleomorphic hepatocellular carcinoma cells, and epithelium lung cancer cells, respectively. We included only TF–gene interactions with a Regulatory Circuit confidence score greater than 0.1. The confidence score indicates the normalized promoter activity level in a given cell type (0: not active, 1: maximally active) (33). For the analysis of mouse pancreatic cell dataset, we obtained the mouse pancreatic TF–gene interactions from CellNet (34). In the construction of the PGRNs, any TF–gene interactions involving unmeasured genes were excluded. In summary, the TF–gene network for human lymphoma, hepatocellular carcinoma cell, and epithelium lung cancer cell lines included 31 392 edges pointing from 515 TFs to 5153 genes, 3868 edges pointing from 413 TFs to 953 genes, and 42 145 edges pointing from 515 TFs to 7125 genes, respectively. The mouse pancreatic PGRN included 2922 edges, involving 95 TFs and 588 genes.

For human PIN, we combined the protein–protein interactions from two databases, namely Enrichr (35) and STRING (36). For mouse pancreatic cells, we obtained mouse (*Mus musculus*) PIN from the STRING database (36). For each TF, we identified its protein partners, defined as proteins that are within a network distance of 2 from the TF in the PIN. When using the STRING database, we included all direct protein partners of TFs, and proteins with a network distance of 2 from TFs with a confidence score reported on STRING larger than 0.5. For human lymphoma, hepatocytes, and lung cancer cells, we identified 11 090 protein partners for a subset of 499 TFs (out of 515 TFs), 10 834 protein partners for a subset of 403 TFs (out of 413 TFs) and 6 175 protein partners for a subset of 504 TFs (out of 515 TFs), respectively. For mouse pancreatic cells, we found 6620 protein partners for a subset of 89 TFs (out of 95 TFs).

Finally, in the construction of the PGRNs, we assigned a directed edge from a TF or from a protein partner of a TF, to every gene regulated by the TF. In summary, the cell type-specific PGRN for human lymphoma cells included 21 488, 617 regulatory edges among 11 161 TFs/proteins and 5153 genes. For hepatocellular carcinoma cells, the PGRN comprised 3726, 393 edges among 10 893 TFs/proteins and 953 genes. For human lung cancer cells, the PGRN comprised 30 656 861 edges among 11 346 TFs/proteins and 7125 genes. For mouse pancreatic cells, the PGRN consisted of 1 417 972 edges among 6661 TFs/proteins and 588 genes. While increasing the size of the PGRN, for example by including lesser confident TF–gene and protein–protein in-

teractions or by including proteins with a network distance from TFs larger than 2, would allow the scoring of a higher number of proteins, such strategy often lowers the accuracy of the protein target predictions.

Gene transcription model. The edges in the PGRN have weights, whose magnitudes represent the strength of the gene regulation and whose signs indicate the direction or the mode of the regulation: positive for gene activation and negative for gene repression. The weights are inferred from the gene expression dataset by adapting a procedure described in our previous method DeltaNet (10,11) (see Figure 1B). The inference of the edge weights is based on an ordinary differential equation (ODE) model of the mRNA production of a gene:

$$\frac{d r_k(t)}{d t} = u_k \prod_{j=1}^n r_j(t)^{a_{kj}} - d_k r_k(t) \quad (1)$$

where $r_k(t)$ is the mRNA concentration of gene k at time t , u_k and d_k denotes the mRNA transcription and degradation rate constants respectively, and a_{kj} denotes the gene regulatory influence (or edge weight) of the j th protein on the k th gene.

While the regulatory edges in the model above usually describe TF–gene interactions, in ProTINA, we further accounted for the (indirect) regulation of a gene by proteins that interact with the TFs. For this purpose, we considered a modified ODE model:

$$\frac{d r_k}{d t} = u_k \left(\prod_{j=1}^{n_{TF}} r_j^{a_{kj}} \prod_{q=1}^{n_P} (r_j r_q)^{b_{kjq}} \right) - d_k r_k \quad (2)$$

where a positive (negative) b_{kjq} describes the activation (repression) of the k th gene by a protein q through its interaction with the TF protein j . The variables n_{TF} and n_P denote the numbers of TFs and their protein partners, respectively. The multiplication of two variables r_j and r_q implies that the regulation of gene k by protein q requires the TF protein j (a non-zero r_j). The model in Equation (2) can be simplified into:

$$\begin{aligned} \frac{d r_k}{d t} &= u_k \left(\prod_{j=1}^{n_{TF}} r_j^{a_{kj} + \sum_q b_{kjq}} \right) \left(\prod_{q=1}^{n_P} r_q^{\sum_j b_{kjq}} \right) - d_k r_k \\ &= u_k \left(\prod_{j=1}^{n_{TF}} r_j^{a_{kj}^*} \right) \left(\prod_{q=1}^{n_P} r_q^{a_{kq}^*} \right) - d_k r_k \\ &= u_k \left(\prod_{j=1}^{n_{TF}+n_P} r_j^{a_{kj}^*} \right) - d_k r_k \end{aligned} \quad (3)$$

where a_{kj}^* denotes the overall regulatory influence of each protein j , including TFs and their protein partners, on the expression of gene k . Note that the model in Equation (3) is mathematically equivalent to that in Equation (1).

By taking the pseudo steady-state assumption (i.e. the synthesis rate of mRNA balances the degradation rate, leading to $dr_k/dt = 0$ in Equation (3)), the inference of edge weights (a_{kj}^*) of the PGRN can be rewritten as the following

linear regression problem (see derivation in (10)):

$$c_{ki} = \sum_{j=1}^{n_{TF}+n_P} a_{kj}^* c_{ji} + p_{ki} \quad (4)$$

where c_{ki} denotes the \log_2 FC expression for gene k in sample i . The variable p_{ki} represents the part of \log_2 FC of gene k expression in sample i that cannot be accounted for by the \log_2 FC of its protein regulators. In other words, p_{ki} indicates the perturbations to the expression of gene k . As detailed below, ProTINA relies on the magnitude and directions of such network perturbations (dysregulations) to identify proteins with altered gene regulatory activity.

The dynamical information contained in time-series gene expression profiles could greatly improve the inference of the edge weights above. But, the pseudo steady-state assumption hinders the application of the linear regression in Equation (4) to time-series data. As previously described in (11), time-series information could be accounted for by adding the following linear constraint on the linear regression problem:

$$s_{ki} = \sum_{j=1}^{n_{TF}+n_P} a_{kj}^* s_{ji} \quad (5)$$

where s_{ki} is the time derivatives (slope) of the \log_2 FC of gene k in sample i . In contrast to Equation (4), Equation (5) was derived without assuming pseudo steady-state, which was necessary to account for the dynamics of gene expressions. The slopes of the \log_2 FC at each sampling time point were computed using a second-order accurate finite difference approximation (37). In summary, the estimation of edge weights in ProTINA involved the following linear regression problem:

$$\mathbf{C}_k = \mathbf{A}_k \mathbf{C}_{R_k} + \mathbf{P}_k \quad (6)$$

$$\mathbf{S}_k = \mathbf{A}_k \mathbf{S}_{R_k} \quad (7)$$

where \mathbf{C}_k and \mathbf{S}_k are the $1 \times m$ vectors of \log_2 FC expressions and time-derivatives of gene k across m samples, the subscript R_k refers to the set of $(n_{TF,k} + n_{P,k})$ protein regulators of gene k in the cell type-specific PGRN, \mathbf{C}_{R_k} and \mathbf{S}_{R_k} denote the $(n_{TF}+n_P, k) \times m$ matrices of \log_2 FCs and their slopes across m samples, \mathbf{A}_k is the $1 \times (n_{TF}+n_P)$ vector of weights for edges in the PGRN pointing to gene k , and \mathbf{P}_k is the $1 \times m$ vector of dysregulation impacts of gene k over m samples.

In ProTINA, the vectors \mathbf{A}_k and \mathbf{P}_k for each gene k in Equations (6) and (7) were estimated by ridge regression. The ridge regression provides a solution to an underdetermined linear regression problem of the standard form: $\mathbf{y} = \mathbf{X}\beta + \varepsilon$, using a penalized least square objective function:

$$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2$$

where λ is a shrinkage parameter for the L^2 -norm penalty. Equations (6) and (7) are rewritten into the standard linear regression problem with $\mathbf{y} = [\mathbf{C}_k \ \mathbf{S}_k]^T$, $\mathbf{X} = [[\mathbf{C}_{R_k} \ \mathbf{S}_{R_k}]^T, [I_m \ \mathbf{0}]^T]$, $\beta = [\mathbf{A}_k \ \mathbf{P}_k]^T$. Before applying the ridge regression, we normalized the vectors of \log_2 FCs and slopes to have a unit

norm. Self-loops were excluded in the regression, and thus the diagonal entries of \mathbf{A}_k were set to 0. In the applications of ProTINA, we employed 10-fold cross validations to determine the optimal λ , one that gives the minimum average prediction error. Here, we used the GLMNET package (38) for both the MATLAB and R versions of ProTINA.

Protein target scoring. In ProTINA, each candidate protein target is assigned a score based on the deviation of the expression of its downstream genes. More specifically, we computed the residuals of the linear regression problem in Equation (6) for each gene k , i.e.

$$\mathbf{r}_k = \mathbf{C}_k - \mathbf{A}_k \mathbf{C}_{R_k} \quad (8)$$

where \mathbf{r}_k is the $1 \times m$ vector of residuals for m samples. For each drug treatment, there often exist multiple gene expression profiles, taken at different time points or different doses. Correspondingly, we evaluated the z -score z_{lk} for each drug treatment l and for each gene k , according to

$$z_{lk} = \frac{\bar{r}_{lk}}{\sigma_k / \sqrt{n_l}} \quad (9)$$

where \bar{r}_{lk} denotes the average residual of gene k among the drug treatment samples, σ_k denotes the sample standard deviation of the residuals in all samples besides the drug treatment, and n_l denotes the number of samples from the drug treatment. A positive (negative) z -score indicates that the expression of gene k in the particular sample was higher (lower) than expected based on the expression of its regulators. The greater the magnitude of the z -score, the more significant is the gene dysregulation.

The target score of a TF or protein for a drug is calculated by combining the z -scores of the target genes in the PGRN, as follows (39):

$$s_{ji} = \frac{\sum_{k=1}^{n_D} w_{kj} z_{ki}}{\sqrt{\sum_{k=1}^{n_D} w_k^2}} \quad (10)$$

where z_{ki} denotes the z -score of gene k and s_{ji} denotes the score of the TF/protein j in the drug treatment sample i . The weighting coefficients w_{kj} are set equal to the edge weights a_{kj} divided by the maximum magnitude of a_{kj} across all j . In other words, the weight w_{kj} reflects the fraction of the regulation of gene k expression that could be attributed to protein j . When w_{kj} (or a_{kj}) and z_{ki} have the same signs, $w_{kj} z_{ki}$ thus takes a positive value. As illustrated in Figure 1C, a positive $w_{kj} z_{ki}$ implies an enhanced regulatory activity of protein j on gene k , since the activation (inhibition) of gene k expression by protein j is stronger in this sample than expected by the PGRN model. In contrast, a negative $w_{kj} z_{ki}$ indicates an attenuation of the regulatory influence of protein j on gene k , since the activation (inhibition) of gene k expression by protein j is weaker than predicted by the PGRN model. Consequently, a highly positive (negative) score s_{ji} is an overall indicator of strongly enhanced (attenuated) regulatory activity of protein j by the drug treatment in sample i (see Figure 1D). The protein targets in each drug treatment sample are ranked in decreasing magnitude of the scores s_{ji} .

DeMAND and differential expression analysis

For DeMAND analysis, we employed the public R sub-routines available from the website: <http://califano.c2b2.columbia.edu/demand>. Following the procedure detailed in the original publication (23), we computed the RMA (Robust Multi-array Average) normalized gene expression values as inputs to the analysis. In DeMAND analysis, we used the same cell type-specific PGRNs as those in ProTINA. For each candidate protein target, DeMAND evaluated the P -value of the deviations in the gene expression relationship between the protein target and each of the genes connected to this protein in the PGRN. The drug targets were ranked in increasing magnitude of the combined P -values.

In differential gene expression (DE) analysis, we calculated the \log_2 FC differential expression of each protein in the PGRN, as described in section Gene expression data above. Here, we used the \log_2 FC values directly as the target scores. Correspondingly, we ranked the candidate protein targets in decreasing magnitude of the \log_2 FC gene expression values.

Performance assessment

For comparing the performance of different methods, we computed the area under the receiver operating characteristic curve (AUROC), i.e. the area under the plot of true positive rate against false positive rate, following the procedure adopted in DREAM challenges (40,41). For each method and each drug treatment, we generated a ranked list of protein targets according to decreasing magnitudes of the protein scores in ProTINA, increasing P -values of network dysregulation from DeMAND, and decreasing magnitudes of \log_2 FC gene expression from DE analysis.

Gene set enrichment analysis

For influenza A virus study, we performed a gene set enrichment analysis (GSEA) of the protein target predictions from ProTINA, DeMAND and DE analysis for the KEGG biological pathways (42), using the R package GAGE (Generally Applicable Gene-set/pathway Enrichment analysis) with Kolmogorov-Smirnov tests (43). In the case of ProTINA and DeMAND, target proteins with zero score were excluded from the GSEA.

Reference protein targets

The reference protein targets of compounds in drug treatment studies were compiled from five different public databases of chemical-protein interactions: DrugBank (44), Therapeutic Target Database (TTD) (45), MATADOR (46), Comparative Toxicogenomics Database (CTD) (47), and STITCH (48). DrugBank and TTD provided information on the mechanism of drug actions as well as the proteins that have physical binding interactions with drugs. Meanwhile, MATADOR, CTD, and STITCH gave interactions between proteins and chemical compounds, curated from text mining and experimental evidences. When retrieving the protein targets of drugs from these databases, we collected proteins that directly bind to the queried drugs. The reference targets for each dataset in this study are provided

in Supplementary material 1. Meanwhile, the reference protein targets for influenza A virus study were obtained from (49), where 1292 host proteins that likely physically bind to viral proteins of influenza type A/WSN/33 in human embryonic kidney cells (HEK293) were identified by whole-genome co-immunoprecipitation assays.

RESULTS

New protein target prediction strategy

ProTINA takes advantage of the availability of comprehensive protein–protein and protein–DNA interaction databases to construct, when possible, a tissue or cell type-specific PGRN. The method considers a PGRN with weighted directed edges (see Figure 1A), describing direct and indirect gene transcriptional regulation by TFs and their protein partners. The edge weights are determined by applying ridge regression using the gene expression data based on a kinetic model of the gene transcriptional process (see Figure 1B and Materials and Methods). Here, a positive weight indicates a gene activation, while a negative weight implies a gene repression. Because of the underlying kinetic model, ProTINA is able to incorporate dynamical gene expression data, a common type of data from drug treatment studies (5,24–26). The scoring of drug targets is based on the enhancement or attenuation of protein–gene regulatory interactions caused by the drug treatment. A drug-induced gene regulatory enhancement occurs when the expression of genes that are positively (negatively) regulated by a candidate target, becomes higher (lower) in drug treated samples than what is predicted by the PGRN model (see Figure 1C). A drug-induced attenuation describes the opposite scenario, where the expression of positively (negatively) regulated genes of a target is lower (higher) than expected from the model. For any given differential gene expression sample, a candidate protein target is scored based on the overall enhancement and/or attenuation of its regulatory influence on the downstream genes (see Figure 1D and Materials and Methods). Thus, a protein target with a more positive (negative) score is considered a more likely target of the drug, in which the drug treatment enhances (attenuates) the gene regulatory activity.

Prediction of known targets of drugs

We tested ProTINA's performance in predicting drug targets using gene expression data from three drug treatment studies employing human and mouse cell lines. The first dataset came from the NCI-DREAM drug synergy study using human diffuse large B cell lymphoma OCI-LY3 (24), the second from the compound genotoxicity study using human liver cancer cells HepG2 (25), and the third from the chromatin-targeting compound study using mouse pancreatic cells (26). We compared ProTINA to the state-of-the-art network-based analytical method DeMAND (23), and to the traditional DE analysis. For the analysis of datasets from human cell lines, we constructed cell-type specific PGRNs by combining human PIN from STRING (36) and Enrichr database (35) and human cell-type specific protein–DNA networks from Regulatory Circuit resource (33). Meanwhile, for the construction of mouse pancreatic

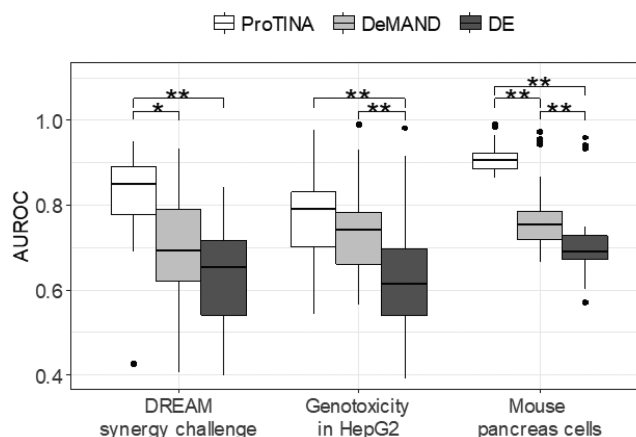


Figure 2. Prediction of known targets of drugs. AUROCs of protein target predictions from ProTINA, DeMAND and DE methods for the NCI-DREAM drug synergy (human B-cell lymphoma), the compound genotoxicity (human HepG2) and the chromatin targeting study (mouse pancreatic cell) datasets (* P -value < 0.01, ** P -value < 0.001 by paired t -test).

cell type-specific PGRN, we used mouse (*Mus musculus*) PIN from STRING (36) and mouse protein–DNA interactions from CellNet (34) (see details in Materials and Methods).

In assessing the performance of ProTINA and the other methods, we compared the ranked list of protein target predictions for each compound with the reference drug targets compiled from the literature (see Materials and Methods and Supplementary material 1). Figure 2 (also see Supplementary Tables S1–S3) summarizes the AUROCs of the target predictions from ProTINA, DeMAND, and DE analysis, showing ProTINA significantly outperforming DeMAND and DE analysis for all three datasets. Here, the drug target predictions from DE analysis had the poorest AUROCs with an overall average below 0.66 (AUROC range: 0.393–0.982). Meanwhile, the target predictions of DeMAND were slightly better than the DE analysis, averaging at 0.74 (AUROC range: 0.405–0.989) for the three datasets. Meanwhile, ProTINA gave the highest average AUROCs among the methods with an average of 0.83 (AUROC range: 0.425–0.991).

Mechanism of action of drugs

Besides high AUROCs, ProTINA also provided accurate and specific indications on the MoA of the compounds. In the NCI-DREAM synergy study, roughly half of the compounds are known to induce DNA damage response, including DNA topoisomerase inhibitors (camptothecin, doxorubicin and etoposide), DNA crosslinker (mitomycin C), oxidative DNA damaging agent (methothrexate), and histone deacetylase (HDAC) inhibitors (trichostatin A). In demonstrating ProTINA's ability to reveal the compound MoA, we focused on the canonical p53 DNA damage response pathway (23), as illustrated in Figure 3. Here, the activation of p53 in response to DNA damage is expected to induce the transcription of Cyclin Dependent Kinase Inhibitor 1A (CDKN1A) and Growth Arrest and DNA Damage Inducible Alpha (GADD45A) (50,51). In turn,

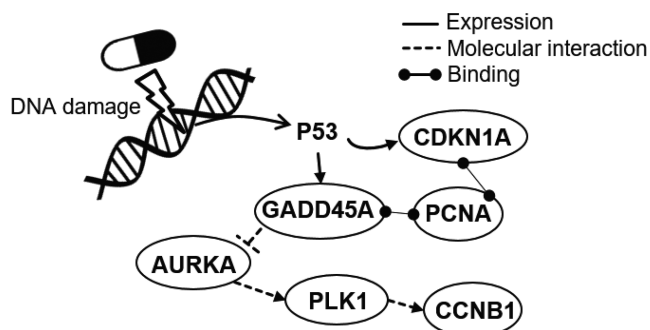


Figure 3. Canonical p53 DNA damage response pathway. In response to DNA damage, GADD45A, CDKN1A, PCNA are activated, while AURKA, CCNB1 and PLK1 proteins are inhibited (23).

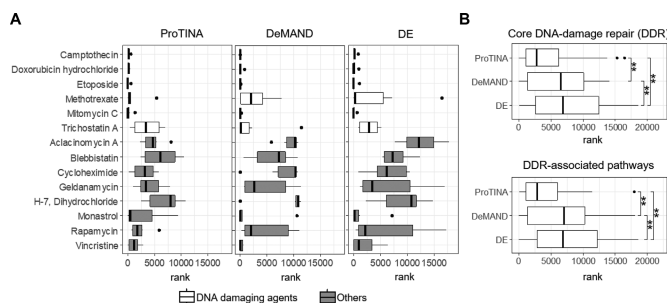


Figure 4. Mechanism of action of compounds based on target predictions by ProTINA. (A) The rank distribution of the canonical p53 DNA damage response proteins in the drug target predictions of ProTINA, DeMAND and DE for the NCI-DREAM drug synergy dataset. (B) The rank distribution of proteins involved in the core DNA-damage repair (DDR) and DDR-associated pathways (56) in the target predictions of ProTINA, DeMAND and DE for the DNA damaging compounds in the NCI-DREAM drug synergy study (***P*-value < 0.001 by Wilcoxon signed rank tests).

CDKN1A and GADD45A—through their interactions with Proliferating Cell Nuclear Antigen (PCNA)—regulate the DNA replication and repair process (52). GADD45A also inhibits the catalytic activity of Aurora Kinase A (AURKA) (53), leading to a lowered activation of Polo-like Kinase 1 (PLK1) and Cyclin B1 (CCNB1) in a phosphorylation cascade (54,55). As shown in Figure 4A, except for trichostatin A, the six proteins in the canonical p53 pathway above were ranked highly by ProTINA among the genotoxic compounds in the study (median rank < 500), consistent with their known MoA. Note that the same six proteins were ranked much lower among the non-DNA damaging compounds (median rank > 500), signifying a high specificity of ProTINA predictions (see also Supplementary Figure S1). Equally important, ProTINA was able to accurately identify the direction of the drug-induced alterations caused by the DNA damaging compounds. The signs of protein target scores from ProTINA indicated drug-induced enhancement (positive scores) of CDKN1A, PCNA and GADD45A, and attenuation (negative scores) of CCNB1, AURKA and PLK1 (see Supplementary Table S4), consistent with the expected response of these proteins to DNA damage in Figure 3.

As illustrated in Figure 4A, DeMAND and DE analysis also performed reasonably well in predicting the com-

pounds' MoA. But, the directions of the perturbations predicted by DE analysis were not always consistent with the expected response to DNA damage (see Supplementary Table S5 and S6). Meanwhile, DeMAND did not provide any information on the directions of the drug perturbations. In addition, the protein target scores of ProTINA provided a clearer demarcation between the genotoxic and the non-genotoxic agents among the compounds in the dataset, than DeMAND and DE analysis (see Supplementary Figure S1). Besides the canonical p53 response pathway, we further looked at the ranking of proteins involved in the overall DNA damage repair (DDR) and its associated pathways (56) (see Supplementary material 2). As depicted in Figure 4B, ProTINA ranked these proteins much higher than DeMAND and DE analysis, with DE performing the poorest among the methods considered.

In comparison to DeMAND and DE analysis, ProTINA was further able to detect a specific MoA of mitomycin C, whose DNA crosslinking activity is expected to prompt a particular DNA repair process called the fanconi anemia pathway (57). The fanconi anemia pathway relies on a specific protein complex to ubiquitinate Fanconi Anemia Group D2 Protein (FANCD2) and Fanconi Anemia Group I Protein (FANCI), as well as two homologous recombination (HR) repair proteins, namely Breast Cancer Type 1 Susceptibility Protein (BRCA1) and RAD51 Recombinase (RAD51) (58). In ProTINA analysis, the average rank of FANCD2, FANCI, BRCA1, and RAD51 was within top 100 for mitomycin C, while the average rank of those proteins was much >100 for the other DNA damaging agents (see Supplementary Table S7). However, the specific activation of the fanconi anemia pathway by mitomycin C was not detected by DeMAND or DE analysis. Thus, ProTINA provided more sensitive and specific indications for the mechanism of action of compounds than DeMAND and DE.

Application of ProTINA for predicting pathogen-host interactions

We applied ProTINA to time-course gene expression profiles of human lung cancer cells (Calu-3) under influenza A virus infection, with the goal of identifying host factors that interact with the viral proteins. The gene expression data came from four studies of influenza A viruses, including A/Netherlands/602/2009 (H1N1), A/CA/04/2009 (H1N1), and A/Vietnam/1203/2004 (H5N1) (27–30). We employed ProTINA to compute the overall protein target scores using the gene expression data of Calu-3 from the four studies above, by averaging the scores from the early phase of the influenza infection between 0 and 12 h. We checked the target predictions of ProTINA against the findings from a genome-wide co-immunoprecipitation analysis of host and viral protein interactions (49). More specifically, the aforementioned study reported 1292 host proteins that co-immunoprecipitated with viral proteins of influenza A/WSN/33 using human embryonic kidney cells (HEK293). Despite the discrepancy in the cell types and influenza viral strains between the co-immunoprecipitation analysis and the gene expression profiling, influenza A viruses share similar features and common protein interac-

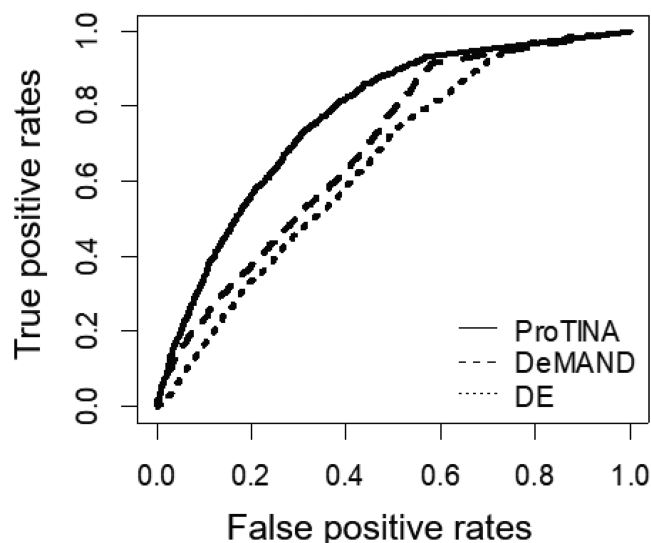


Figure 5. Prediction of targets of influenza A virus. The receiver operative characteristic curves give the true positive rate versus the false positive rate relationship of the protein target predictions from ProTINA, DeMAND and DE against proteins that co-immunoprecipitate with influenza A viral proteins. The AUROCs for ProTINA, DeMAND and DE analysis are 0.77, 0.69 and 0.65, respectively.

tions (59,60). Besides ProTINA, we also evaluated the accuracy of viral target predictions from DeMAND and DE for the same dataset.

Figure 5 gives the receiver operating characteristic (ROC) curves of the target predictions from ProTINA, DeMAND and DE analysis. ProTINA outperformed the two other methods, providing the highest AUROC (ProTINA: 0.76 versus demand: 0.69 and DE: 0.65). We further performed a gene set enrichment analysis (GSEA) for the target predictions from each of the methods (see Materials and Methods) to elucidate the key pathways involved in the viral infection and the accompanying host response. The results of the GSEA are summarized in Figure 6. Both DeMAND and DE target predictions were enriched for only a few pathways (q -value < 0.01), while ProTINA prediction had a much higher number of overrepresented pathways.

The common enriched pathways among ProTINA, DeMAND and DE (top of Figure 6) included known mechanisms related to viral entry, replication and assembly, including endocytosis (61), protein processing in endoplasmic reticulum (62), ubiquitin mediated proteolysis (63,64) and RIG-I-like receptor signaling pathway (65,66). Both ProTINA and DE analysis indicated the modulation of host cell cycle (67), mRNA surveillance (68) and DNA damage response (69). Only ProTINA prediction was significantly enriched for focal adhesion and actin cytoskeleton, which have been shown to regulate influenza virus entry at the early stage of infection (70). In addition, ProTINA target predictions were also enriched for a broad set of host response pathways to viral infection, including host defense mechanism (e.g. T- and B-cell receptor pathways, phagocytosis, leukocyte migration, chemokine signaling pathways), DNA damage repair (e.g. nucleotide excision repair, p53 signaling pathway, homologous recombination) and apop-

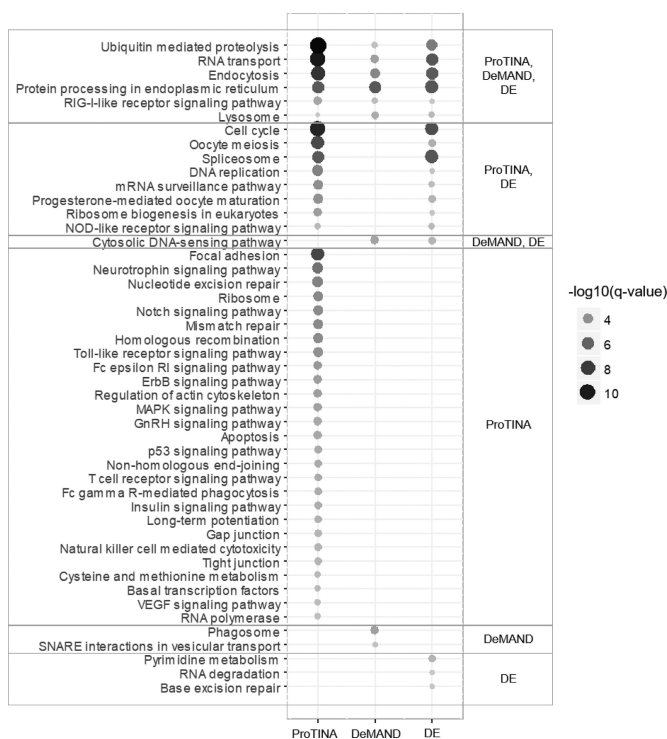


Figure 6. Gene set enrichment analysis for KEGG pathways for the influenza A protein target predictions from ProTINA, DeMAND and DE. The size of the circles corresponds to $-\log_{10}$ scale of the q -values. Only pathways with q -value < 0.01 are shown.

toxis. As several influenza proteins are known to interfere with interferon production (which in turn regulates several cytokines) (65,66), these findings suggest that, overall, ProTINA provided a broader picture of the early events in the influenza A viral infection, than DeMAND and DE analysis.

DISCUSSION

ProTINA is a novel and highly effective network-based analytical method for inferring the protein targets of compounds from gene expression profiling data. ProTINA combines the information of TF–gene and protein–protein interactions and data of differential gene expressions to create a tissue or cell type-specific PGRN model. Similar to network-based analysis methods such as NIR (7), MNI (8), SSEM (9) and DeltaNet (10), ProTINA uses a dynamic mechanistic model of the gene transcriptional process to compute deviations in the differential gene expression profiles that are induced by drug treatments. However, as mentioned earlier, the expression of the targets of a drug is often unaffected by the drug treatment (3). For this reason and as illustrated in Figure 1C and D, ProTINA further transforms the deviations in the differential gene expression into alterations in the protein–gene regulatory edges in the PGRN model. Finally, the target scoring is based on edgetic perturbations of the PGRN, specifically enhancement or attenuation of gene regulatory interactions, caused by the compound.

Like ProTINA, the state-of-the-art method DeMAND also relies on gene transcriptional dysregulations to score drug targets. But, DeMAND does not consider the mode nor the dynamics of the gene regulations, and is unable to predict the direction of the drug-induced dysregulations. DeMAND calculates protein dysregulation scores (P -values) for a given gene regulatory network, by statistical comparison between samples from drug treatment and from control experiments. Consequently, DeMAND requires only few samples to generate its prediction (provided that the network can be defined *a priori*). On the other hand, ProTINA makes use of the available differential gene expression profiles from a study or a cell line (i.e. not only from the specific drug treatment), to assign the edge weights of the PGRN by ridge regression. Importantly, in the regression analysis, the PGRN model used in ProTINA accounts for the network perturbations. The ability of ProTINA to incorporate data from other drug treatments or conditions in the scoring of protein targets makes this method particularly suited to take advantage of the extensive and still growing number of gene transcriptional profiles from online databases, such as GEO. As demonstrated in the applications to three benchmark drug treatment datasets using human and mouse cell lines, ProTINA significantly outperformed DeMAND and the standard DE analysis. The target predictions of ProTINA also provide indications for the MoA of compounds, including the directions of the network perturbations, with high sensitivity and specificity.

Besides its intended use to predict targets of compounds, we also demonstrated that the analysis of network perturbations using ProTINA could provide insights into the mechanism of diseases. In the application to gene expression profiles of Calu-3 cells from influenza A infection studies, ProTINA again outperformed DeMAND and DE analysis in identifying host factors that bind with viral proteins. Furthermore, the GSEA of ProTINA target predictions revealed the spectrum of cellular processes involved in the early phase of influenza A infection, including pathways involved in viral entry, replication and assembly, and those related to cellular response to viral infection. Among the pathways with the highest significance (lowest q -value) was focal adhesion, which has been shown to regulate influenza viral entry as well as viral replication (70). Meanwhile, the target predictions of DeMAND and DE analysis had fewer enriched pathways, and thus were less informative than the target analysis by ProTINA.

The PGRN model (see Equation (1)) belongs to a class of modeling framework called Biochemical Systems Theory, specifically the S-systems model (71). In addition to gene regulatory networks, S-system modeling have also been used to describe other cellular networks, including signal transduction pathways and metabolic reaction networks (72). Therefore, the principle used in ProTINA could be readily adapted to infer perturbations in cellular signalling or metabolic networks, for example from proteomic and metabolomics profiles, respectively. Besides PGRNs and gene transcriptional profiles, we have not applied ProTINA to analyze other types of cellular networks and data, as such an application was beyond the scope of our work.

ProTINA requires a cell type- or tissue-specific PGRN as an input, which may hinder its application to analyze

data from lesser studied organisms. In the case studies, we leveraged on the extensive online databases of protein–protein interactions and TF–gene networks to manually curate PGRNs for human and mouse cells (33,34,36). Alternatively, provided that a large dataset of gene expression profiles are available for the cell of interest, the PGRN could be inferred using existing network inference methods (73,74). Another potential limitation in applying ProTINA is the requirement for differential gene expression data for inferring the edge weights of the PGRN. While the minimum number for implementing ridge regression with a 3-fold cross validation (lowest fold in GLMNET) is three, the accuracy of the weights and thus the target predictions from ProTINA would generally deteriorate with lower sample sizes. Nevertheless, ProTINA was still able to provide reasonably accurate predictions using a total of 18 samples in the influenza A virus case study above.

The performance of ProTINA, like any other network-based analytical methods, depends on the fidelity of the network used in the analysis. Uncertainty in the PGRN model, both in the structure and the edge weights, is expected to negatively affect the accuracy of the target prediction. Here, structural uncertainty is associated with the reliability of the information used to construct the PGRN, which in our study, comes from online databases of PIN and TF–gene networks. On the other hand, the uncertainty in the edge weights is associated with multiple factors, including the information content of the gene transcriptional profiles and the mathematical formulation used for the weight inference. The information content of the gene expression data is in turn related to measurement uncertainty and richness in the experimental perturbations. Keeping the same number of treatments, datasets with more replicates and less correlated gene expression profiles (i.e. the treatments induce more distinct perturbations to the network), would have a higher degree of information. Meanwhile, we have previously shown that the validity of the model assumption (e.g. pseudo steady-state condition) has an effect on the accuracy of the inferred weights and thus the target prediction accuracy (10). While we have circumvented the issue arising from the violation of the pseudo steady-state assumption in ProTINA (see Equation (7)), (in)validating all model assumptions may be difficult, if not impossible, in practice. A common strategy, as implemented in this study, is to test the performance of the method against benchmark datasets (13). The results of applying ProTINA to drug treatment and influenza A viral infection datasets give confidence to the suitability of the mathematical formulation used in this work.

AVAILABILITY

MATLAB and R versions of ProTINA can be downloaded from Github repository (<https://github.com/CABSEL/ProTINA>).

SUPPLEMENTARY DATA

Supplementary Data are available at NAR online.

ACKNOWLEDGEMENT

We would like to thank Ziyi Hua for her assistance in preparing the R codes of ProTINA.

FUNDING

ETH Research Grant (project title: “Network Perturbation Discovery under Uncertainty”). Funding for open access charge: ETH Zurich.

Conflict of interest statement. None declared.

REFERENCES

- Imming,P., Sinning,C. and Meyer,A. (2006) Drugs, their targets and the nature and number of drug targets. *Nat. Rev. Drug Discov.*, **5**, 821–834.
- Schenone,M., Dančik,V., Wagner,B.K. and Clemons,P.A. (2013) Target identification and mechanism of action in chemical biology and drug discovery. *Nat. Chem. Biol.*, **9**, 232–240.
- Isik,Z., Baldow,C., Cannistraci,C.V. and Schroeder,M. (2015) Drug target prioritization by perturbed gene expression and network information. *Sci. Rep.*, **5**, 17417.
- Chua,H.N. and Roth,F.P. (2011) Discovering the targets of drugs via computational systems biology. *J. Biol. Chem.*, **286**, 23653–23658.
- Lamb,J. (2007) The connectivity map: A new tool for biomedical research. *Nat. Rev. Cancer*, **7**, 54–60.
- Iorio,F., Rittman,T., Ge,H., Menden,M. and Saez-Rodriguez,J. (2013) Transcriptional data: a new gateway to drug repositioning? *Drug Discov. Today*, **18**, 350–357.
- Gardner,T.S. (2003) Inferring genetic networks and identifying compound mode of action via expression profiling. *Science*, **301**, 102–105.
- di Bernardo,D., Thompson,M.J., Gardner,T.S., Chobot,S.E., Eastwood,E.L., Wojtovich,A.P., Elliott,S.J., Schaus,S.E. and Collins,J.J. (2005) Chemogenomic profiling on a genome-wide scale using reverse-engineered gene networks. *Nat. Biotechnol.*, **23**, 377–383.
- Cosgrove,E.J., Zhou,Y., Gardner,T.S. and Kolaczyk,E.D. (2008) Predicting gene targets of perturbations via network-based filtering of mRNA expression compendia. *Bioinformatics*, **24**, 2482–2490.
- Noh,H. and Gunawan,R. (2016) Inferring gene targets of drugs and chemical compounds from gene expression profiles. *Bioinformatics*, **32**, 2120–2127.
- Noh,H., Ziyi,H. and Gunawan,R. (2016) Inferring causal gene targets from time course expression data. *IFAC-PapersOnLine*, **49**, 350–356.
- Brock,A., Krause,S., Li,H., Kowalski,M., Goldberg,M.S., Collins,J.J. and Ingber,D.E. (2014) Silencing HoxA1 by intraductal injection of siRNA lipidoid nanoparticles prevents mammary tumor progression in mice. *Sci. Transl. Med.*, **6**, 217ra2.
- Marbach,D., Costello,J.C., Küffner,R., Vega,N.M., Prill,R.J., Camacho,D.M., Allison,K.R., Kellis,M., Collins,J.J., Aderhold,A. *et al.* (2012) Wisdom of crowds for robust gene network inference. *Nat. Methods*, **9**, 796–804.
- Siegenthaler,C. and Gunawan,R. (2014) Assessment of network inference methods: How to cope with an underdetermined problem. *PLoS One*, **9**, e90481.
- Ud-Dean,S.M.M. and Gunawan,R. (2014) Ensemble inference and inferability of gene regulatory networks. *PLoS One*, **9**, e103812.
- Lachmann,A. and Maayan,A. (2009) KEA: Kinase enrichment analysis. *Bioinformatics*, **25**, 684–686.
- Lefebvre,C., Rajbhandari,P., Alvarez,M.J., Bandaru,P., Lim,W.K., Sato,M., Wang,K., Sumazin,P., Kustagi,M., Bisikirska,B.C. *et al.* (2010) A human B-cell interactome identifies MYB and FOXM1 as master regulators of proliferation in germinal centers. *Mol. Syst. Biol.*, **6**, 377.
- Chen,E.Y., Xu,H., Gordonov,S., Lim,M.P., Perkins,M.H. and Ma'ayan,A. (2017) Expression2Kinases: mRNA profiling linked to multiple upstream regulatory layers. *Bioinformatics*, **28**, 105–111.
- Emig,D., Ivliev,A., Pustovalova,O., Lancashire,L., Bureeva,S., Nikolsky,Y. and Bessarabova,M. (2013) Drug target prediction and repositioning using an integrated network-based approach. *PLoS One*, **8**, e60618.
- Laenen,G., Thorrez,L., Börnigen,D. and Moreau,Y. (2013) Finding the targets of a drug by integration of gene expression data with a protein interaction network. *Mol. Biosyst.*, **9**, 1676–1685.
- Chindelevitch,L., Ziemek,D., Enayetallah,A., Randhawa,R., Sidders,B., Brockel,C. and Huang,E.S. (2012) Causal reasoning on biological networks: Interpreting transcriptional changes. *Bioinformatics*, **28**, 1114–1121.
- Martin,F., Thomson,T.M., Sewer,A., Drubin,D.A., Mathis,C., Weisensee,D., Pratt,D., Hoeng,J. and Peitsch,M.C. (2012) Assessment of network perturbation amplitudes by applying high-throughput data to causal biological networks. *BMC Syst. Biol.*, **6**, 54.
- Woo,J.H., Shimoni,Y., Yang,W.S., Subramaniam,P., Iyer,A., Nicoletti,P., Rodríguez Martínez,M., López,G., Mattioli,M., Realubit,R. *et al.* (2015) Elucidating compound mechanism of action by network perturbation analysis. *Cell*, **162**, 441–451.
- Bansal,M., Yang,J., Karan,C., Menden,M.P., Costello,J.C., Tang,H., Xiao,G., Li,Y., Allen,J., Zhong,R. *et al.* (2014) A community computational challenge to predict the activity of pairs of compounds. *Nat. Biotechnol.*, **32**, 1213–1222.
- Magkoufopoulou,C., Claessen,S.M.H., Tsamou,M., Jennen,D.G.J., Kleinjans,J.C.S. and Van delft,J.H.M. (2012) A transcriptomics-based in vitro assay for predicting chemical genotoxicity in vivo. *Carcinogenesis*, **33**, 1421–1429.
- Kubicek,S., Gilbert,J.C., Fomina-yadlin,D., Gitlin,A.D. and Yuan,Y. (2012) Chromatin-targeting small molecules cause class-specific transcriptional changes in pancreatic endocrine cells. *Proc. Natl. Acad. Sci. U.S.A.*, **109**, 5364–5369.
- McDermott,J.E., Shankaran,H., Eisfeld,A.J., Belisle,S.E., Neuman,G., Li,C., McWeeny,S., Sabourin,C., Kawaoka,Y., Katze,M.G. *et al.* (2011) Conserved host response to highly pathogenic avian influenza virus infection in human cell culture, mouse and macaque model systems. *BMC Syst. Biol.*, **5**, 190.
- Li,C., Bankhead,A., Eisfeld,A.J., Hatta,Y., Jeng,S., Chang,J.H., Aicher,L.D., Proll,S., Ellis,A.L., Law,G.L. *et al.* (2011) Host regulatory network response to infection with highly pathogenic H5N1 avian influenza virus. *J. Virol.*, **85**, 10955–10967.
- Mitchell,H.D., Eisfeld,A.J., Sims,A.C., McDermott,J.E., Matzke,M.M., Webb-Robertson,B.J.M., Tilton,S.C., Tchitchek,N., Josset,L., Li,C. *et al.* (2013) A network integration approach to predict conserved regulators related to pathogenicity of influenza and SARS-CoV respiratory viruses. *PLoS One*, **8**, e69374.
- Menachery,V.D., Eisfeld,A.J., Schäfer,A., Josset,L., Sims,A.C., Proll,S., Fan,S., Li,C., Neumann,G., Tilton,S.C. *et al.* (2014) Pathogenic influenza viruses and coronaviruses utilize similar and contrasting approaches to control interferon-stimulated gene responses. *MBio*, **5**, e01174-14.
- Barrett,T., Wilhite,S.E., Ledoux,P., Evangelista,C., Kim,I.F., Tomashevsky,M., Marshall,K.A., Phillippy,K.H., Sherman,P.M., Holko,M. *et al.* (2013) NCBI GEO: archive for functional genomics data sets - Update. *Nucleic Acids Res.*, **41**, 991–995.
- Gentleman,R., Carey,V., Bates,D., Bolstad,B., Dettling,M., Dudoit,S., Ellis,B., Gautier,L., Ge,Y., Gentry,J. *et al.* (2004) Bioconductor: open software development for computational biology and bioinformatics. *Genome Biol.*, **5**, R80.
- Marbach,D., Lamparter,D., Quon,G., Kellis,M., Kutalik,Z. and Bergmann,S. (2016) Tissue-specific regulatory circuits reveal variable modular perturbations across complex diseases. *Nat. Methods*, **13**, 366–370.
- Cahan,P., Li,H., Morris,S.A., Lummertz Da Rocha,E., Daley,G.Q. and Collins,J.J. (2014) CellNet: Network biology applied to stem cell engineering. *Cell*, **158**, 903–915.
- Kuleshov,M.V., Jones,M.R., Rouillard,A.D., Fernandez,N.F., Duan,Q., Wang,Z., Koplev,S., Jenkins,S.L., Jagodnik,K.M., Lachmann,A. *et al.* (2016) Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.*, **44**, W90–W97.
- Szklarczyk,D., Franceschini,A., Wyder,S., Forslund,K., Heller,D., Huerta-Cepas,J., Simonovic,M., Roth,A., Santos,A., Tsafou,K.P. *et al.* (2015) STRING v10: protein–protein interaction networks, integrated over the tree of life. *Nucleic Acids Res.*, **43**, D447–D452.
- Lynch,D.R., (2005) *Finite Difference Calculus*, Numerical Partial Differential Equations for Environmental Scientists and Engineers: A first Practical Course. pp. 11–19.

38. Friedman, J., Hastie, T. and Tibshirani, R. (2010) Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.*, **33**, 1–22.
39. Whitlock, M.C. (2005) Combining probability from independent tests: The weighted Z-method is superior to Fisher's approach. *J. Evol. Biol.*, **18**, 1368–1373.
40. Stolovitzky, G., Prill, R.J. and Califano, A. (2009) Lessons from the DREAM2 challenges: a community effort to assess biological network inference. *Ann. N. Y. Acad. Sci.*, **1158**, 159–195.
41. Prill, R.J., Marbach, D., Saez-Rodriguez, J., Sorger, P.K., Alexopoulos, L.G., Xue, X., Clarke, N.D., Altan-Bonnet, G. and Stolovitzky, G. (2010) Towards a rigorous assessment of systems biology models: The DREAM3 challenges. *PLoS One*, **5**, e9202.
42. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. and Morishima, K. (2017) KEGG: New perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.*, **45**, D353–D361.
43. Luo, W., Friedman, M.S., Shedden, K., Hankenson, K.D. and Woolf, P.J. (2009) GAGE: generally applicable gene set enrichment for pathway analysis. *BMC Bioinformatics*, **10**, 161.
44. Law, V., Knox, C., Djoumbou, Y., Jewison, T., Guo, A.C., Liu, Y., MacIejewski, A., Arndt, D., Wilson, M., Neveu, V. et al. (2014) DrugBank 4.0: Shedding new light on drug metabolism. *Nucleic Acids Res.*, **42**, 1091–1097.
45. Zhu, F., Shi, Z., Qin, C., Tao, L., Liu, X., Xu, F., Zhang, L., Song, Y., Liu, X., Zhang, J. et al. (2012) Therapeutic target database update 2012: a resource for facilitating target-oriented drug discovery. *Nucleic Acids Res.*, **40**, 1128–1136.
46. Günther, S., Kuhn, M., Dunkel, M., Campillos, M., Senger, C., Petsalaki, E., Ahmed, J., Urdiales, E.G., Gewiss, A., Jensen, L.J. et al. (2008) SuperTarget and Matador: resources for exploring drug-target relationships. *Nucleic Acids Res.*, **36**, 919–922.
47. Davis, A.P., Grondin, C.J., Johnson, R.J., Sciaky, D., King, B.L., McMorran, R., Wiegiers, J., Wiegiers, T.C. and Mattingly, C.J. (2017) The comparative toxicogenomics database: update 2017. *Nucleic Acids Res.*, **45**, D972–D978.
48. Szklarczyk, D., Santos, A., Von Mering, C., Jensen, L.J., Bork, P. and Kuhn, M. (2016) STITCH 5: augmenting protein-chemical interaction networks with tissue and affinity data. *Nucleic Acids Res.*, **44**, D380–D384.
49. Watanabe, T., Kawakami, E., Shoemaker, J.E., Lopes, T.J.S., Matsuoka, Y., Tomita, Y., Kozuka-Hata, H., Gorai, T., Kuwahara, T., Takeda, E. et al. (2014) Influenza virus-host interactome screen as a platform for antiviral drug development. *Cell Host Microbe*, **16**, 795–805.
50. Cazzalini, O., Scovassi, A.I., Savio, M., Stivala, L.A. and Prosperi, E. (2010) Multiple roles of the cell cycle inhibitor p21CDKN1A in the DNA damage response. *Mutat. Res.*, **704**, 12–20.
51. Zhan, Q. (2005) Gadd45a, a p53- and BRCA1-regulated stress protein, in cellular response to DNA damage. *Mutat. Res. - Fundam. Mol. Mech. Mutagen.*, **569**, 133–143.
52. Kelman, Z. (1997) PCNA: structure, functions and interactions. *Oncogene*, **14**, 629–640.
53. Shao, S., Wang, Y., Jin, S., Song, Y., Wang, X., Fan, W., Zhao, Z., Fu, M., Tong, T., Dong, L. et al. (2006) Gadd45a interacts with aurora-A and inhibits its kinase activity. *J. Biol. Chem.* **281**, 28943–28950.
54. Macůrek, L., Lindqvist, A., Lim, D., Lampson, M.A., Klompaker, R., Freire, R., Clouin, C., Taylor, S.S., Yaffe, M.B. and Medema, R.H. (2008) Polo-like kinase-1 is activated by aurora A to promote checkpoint recovery. *Nature*, **455**, 119–123.
55. Toyoshima-Morimoto, F., Taniguchi, E., Shinya, N., Iwamatsu, A. and Nishida, E. (2001) Polo-like kinase 1 phosphorylates cyclin B1 and targets it to the nucleus during prophase. *Nature*, **410**, 215–220.
56. Pearl, L.H., Schierz, A.C., Ward, S.E., Al-Lazikani, B. and Pearl, F.M.G. (2015) Therapeutic opportunities within the DNA damage response. *Nat. Rev. Cancer*, **15**, 166–180.
57. Deans, A.J. and West, S.C. (2013) DNA interstrand crosslink repair and cancer. *Nat rev Cancer*, **11**, 467–480.
58. Andreassen, P. and Ren, K. (2009) Fanconi anaemia proteins, DNA interstrand crosslink repair pathways, and cancer therapy. *Curr. Cancer Crug Targets*, **9**, 101–117.
59. Shaw, M.L., Stone, K.L., Colangelo, C.M., Gulcicek, E.E. and Palese, P. (2008) Cellular proteins in influenza virus particles. *PLoS Pathog.*, **4**, e1000085.
60. de Chasse, B., Meyniel-Schicklin, L., Aublin-Gex, A., Navratil, V., Chantier, T., André, P. and Lotteau, V. (2013) Structure homology and interaction redundancy for discovering virus–host protein interactions. *EMBO Rep.*, **14**, 938–944.
61. Matsuoka, Y., Matsumae, H., Katoh, M., Eisfeld, A.J., Neumann, G., Hase, T., Ghosh, S., Shoemaker, J.E., Lopes, T.J., Watanabe, T. et al. (2013) A comprehensive map of the influenza A virus replication cycle. *BMC Syst. Biol.*, **7**, 97.
62. Braakman, I., Hoover-Litty, H., Wagner, K.R. and Helenius, A. (1991) Folding of influenza hemagglutinin in the endoplasmic reticulum. *J. Cell Biol.*, **114**, 401–411.
63. Rodriguez, A., Perez-Gonzalez, A. and Nieto, A. (2007) Influenza virus infection causes specific degradation of the largest subunit of cellular RNA polymerase II. *J. Virol.*, **81**, 5315–5324.
64. Rudnicka, A. and Yamauchi, Y. (2016) Ubiquitin in influenza virus entry and innate immunity. *Viruses*, **8**, 293.
65. Hale, B.G., Randall, R.E., Ortin, J. and Jackson, D. (2008) The multifunctional NS1 protein of influenza A viruses. *J. Gen. Virol.*, **89**, 2359–2376.
66. Varga, Z.T., Grant, A., Manicassamy, B. and Palese, P. (2012) Influenza virus protein PB1-F2 inhibits the induction of type I interferon by binding to MAVS and decreasing mitochondrial membrane potential. *J. Virol.*, **86**, 8359–8366.
67. Shoemaker, J.E., Fukuyama, S., Eisfeld, A.J., Muramoto, Y., Watanabe, S., Watanabe, T., Matsuoka, Y., Kitano, H. and Kawaoka, Y. (2012) Integrated network analysis reveals a novel role for the cell cycle in 2009 pandemic influenza virus-induced inflammation in macaque lungs. *BMC Syst. Biol.*, **6**, 117.
68. Cho, H., Ahn, S.H., Kim, K.M. and Kim, Y.K. (2013) Non-structural protein 1 of influenza viruses inhibits rapid mRNA degradation mediated by double-stranded RNA-binding protein, staufen1. *FEBS Lett.*, **587**, 2118–2124.
69. Li, N., Parrish, M., Chan, T.K., Yin, L., Rai, P., Yoshiyuki, Y., Abolhassani, N., Tan, K.B., Kiraly, O., Chow, V.T.K. et al. (2015) Influenza infection induces host DNA damage and dynamic DNA damage responses during tissue regeneration. *Cell. Mol. Life Sci.*, **72**, 2973–2988.
70. Elbahesh, H., Cline, T., Baranovich, T., Govorkova, E.A., Schultz-Cherry, S. and Russell, C.J. (2014) Novel roles of focal adhesion kinase in cytoplasmic entry and replication of influenza A viruses. *J. Virol.* **88**, 6714–6728.
71. Voit, E.O. (2000) *Computational Analysis of Biochemical Systems: A Practical Guide for Biochemists and Molecular Biologists*, Cambridge University Press, Cambridge, UK.
72. Chou, I.C. and Voit, E.O. (2009) Recent developments in parameter estimation and structure identification of biochemical and genomic systems. *Math. Biosci.*, **219**, 57–83.
73. Margolin, A.A., Nemenman, I., Basso, K., Wiggins, C., Stolovitzky, G., Favera, R.D. and Califano, A. (2006) ARACNE: An algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinformatics*, **7**, S7.
74. Wang, K., Saito, M., Bisikirska, B.C., Alvarez, M.J., Lim, W.K., Rajbhandari, P., Shen, Q., Nemenman, I., Basso, K., Margolin, A.A. et al. (2009) Genome-wide identification of post-translational modulators of transcription factor activity in human B cells. *Nat. Biotechnol.*, **27**, 829–837.