

Identifying proteins controlling key disease signaling pathways

Anthony Gitter^{1,†} and Ziv Bar-Joseph^{2,*}

¹Computer Science Department and ²Lane Center for Computational Biology, School of Computer Science, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213, USA

ABSTRACT

Motivation: Several types of studies, including genome-wide association studies and RNA interference screens, strive to link genes to diseases. Although these approaches have had some success, genetic variants are often only present in a small subset of the population, and screens are noisy with low overlap between experiments in different labs. Neither provides a mechanistic model explaining how identified genes impact the disease of interest or the dynamics of the pathways those genes regulate. Such mechanistic models could be used to accurately predict downstream effects of knocking down pathway members and allow comprehensive exploration of the effects of targeting pairs or higher-order combinations of genes.

Results: We developed methods to model the activation of signaling and dynamic regulatory networks involved in disease progression. Our model, SDREM, integrates static and time series data to link proteins and the pathways they regulate in these networks. SDREM uses prior information about proteins' likelihood of involvement in a disease (e.g. from screens) to improve the quality of the predicted signaling pathways. We used our algorithms to study the human immune response to H1N1 influenza infection. The resulting networks correctly identified many of the known pathways and transcriptional regulators of this disease. Furthermore, they accurately predict RNA interference effects and can be used to infer genetic interactions, greatly improving over other methods suggested for this task. Applying our method to the more pathogenic H5N1 influenza allowed us to identify several strain-specific targets of this infection.

Availability: SDREM is available from <http://sb.cs.cmu.edu/sdrem>

Contact: zivbj@cs.cmu.edu

Supplementary information: Supplementary data are available at *Bioinformatics* online.

1 INTRODUCTION

A wide variety of experimental and computational approaches have been used over the past few years to screen for genes that play important roles in human disease. These include RNA interference (RNAi) screens (Mohr *et al.*, 2010), which knock down human genes to quantify the phenotypic effects in diseases such as HIV (Brass *et al.*, 2008) and influenza (Brass *et al.*, 2009; Karlas *et al.*, 2010; König *et al.*, 2010) infection, and genome-wide association studies (GWAS), a powerful approach for uncovering strong connections between genetic variations such as single-nucleotide polymorphisms (SNPs) and disease traits (Altshuler *et al.*, 2008).

Although these and similar sequencing-based methods have been widely applied and were able to identify many relevant genes, they also suffer from important drawbacks that limit

their effectiveness. RNAi screen hits are often not reproducible among different labs and experimental settings. For example, the overlap among hits in HIV screens is low. Three genome-wide screens identified 842 genes that impact HIV replication, but only three genes were common to all screens (Bushman *et al.*, 2009). Similar low overlap is observed for other diseases as well (Sertiz and Shaw, 2011). Interpretability of screen hits remains a challenge, as the screens cannot explain why a gene impacts a disease. Similarly, interpretation of GWAS hits is not straightforward as they oftentimes lie outside of coding regions of the genome (Schaub *et al.*, 2012). Even those variants lying within a coding region only explain a small fraction of affected individuals for several diseases and conditions (Maher, 2008).

These observations motivated methods that integrate GWAS or perturbation information with functional annotations and other types of genomic data (most notably gene expression and protein interaction data). For example, pathway-based GWAS (Wang *et al.*, 2010) maps SNP *P*-values to proximal genes and links pathways to diseases by testing whether any predefined pathways are enriched for genes that flank significant SNPs. However, this and similar strategies are dependent on the quality of the annotated pathways preventing their use for important diseases that are poorly represented in public pathway databases (Wang *et al.*, 2010).

A variety of gene prioritization algorithms, which have been extensively reviewed (Moreau and Tranchevent, 2012; Piro and Di Cunto, 2012) and benchmarked (Börnigen *et al.*, 2012), use known disease genes from literature curation, GWAS or other sources as seeds to search for similar genes that are likely to be related to the same disease. The evidence used to define gene similarity includes text mining, pathway membership, functional annotations, protein properties, sequence, co-expression and proximity in protein–protein interaction (PPI) networks. However, these methods do not attempt to mechanistically model the pathways that are activated during disease progression and response and thus do not provide a model to explain the various observations. Instead, they use a 'black box' approach; they take seed genes as input and produce a ranked list of candidate genes without explaining the predicted relationship to the disease, which may impair further analysis (Moreau and Tranchevent, 2012).

Other algorithms use perturbation and genetic data to link proteins to diseases, while at the same time using other genomic datasets and interaction networks to suggest how these proteins may be involved in disease response (Huang and Fraenkel, 2009; Kim *et al.*, 2011; Ourfali *et al.*, 2007; Tuncbag *et al.*, 2013; Yeang *et al.*, 2004; Yeager-Lotem *et al.*, 2009). However, these methods have only used static models so far. In such models, the phenotypic outcomes of the network perturbations (or genetic variants) are typically differentially expressed genes. The algorithms then

*To whom correspondence should be addressed.

[†]Present address: Microsoft Research, 1 Memorial Drive, Cambridge, MA 02142, USA.

connect possible disease sources (for example, genetic variants, perturbed genes, receptors or proteins that directly interact with a pathogen) with the differentially expressed genes and use the resulting source–target pathways to explain the role various proteins play in the disease response.

Static models ignore the dynamic nature of responses and the temporal changes in active genes and pathways. Indeed, many transcriptional studies of disease highlight the importance of the temporal aspects, for example, the rapid immune response following infection (Li *et al.*, 2011; Shapira *et al.*, 2009) and longitudinal clinical studies (Desai *et al.*, 2011). Static modeling of dynamic interactions aggregates several different pathways at once, which leads to missing key factors that are only active at specific times during response (Ernst *et al.*, 2007). In addition, many static algorithms do not attempt to infer the directionality of the edges on the pathways despite the importance of edge orientation in correctly representing signaling paths (Gitter *et al.*, 2011).

To address these issues, we developed the Signaling and Dynamic Regulatory Events Miner (SDREM) (Gitter *et al.*, 2013). SDREM combines an input–output hidden Markov model (IOHMM) (Bengio and Frasconi, 1995) with a combinatorial optimization algorithm to reconstruct dynamic models of the signaling and regulatory networks using high-throughput data including PPI, transcription factor (TF) binding and time series gene expression data. The signaling pathways inferred by SDREM are directed cascades that connect sources and target TFs. These TFs are in turn used as part of the reconstructed dynamic regulatory network to explain the temporal gene expression. SDREM was developed and tested on yeast and was shown to accurately reconstruct yeast stress response networks leading to new insights into the specific pathways controlling such responses (Gitter *et al.*, 2013).

Although SDREM was successfully applied to model organisms, the application to human disease remained a challenge. The larger scale of the human interaction network makes search and inference much more complicated. In addition, several data sources available for studying diseases (including RNAi screens and GWAS data) were not used by the original SDREM algorithm. Finally, human networks are more complex and involve extensive redundant parallel pathways (Logue and Morrison, 2012), which means that new methods are required to assess the ability of proteins and combinations of proteins to affect disease progression.

Here we present new computational methods that extend SDREM by designing a new target function that uses prior node (protein) information, developing new optimization methods for searching large networks and creating new gene ranking metrics based on target connectivity as a post-processing step. Combined, these extensions allowed us to apply SDREM to study human response to influenza A viral infection. The resulting networks identified several key proteins and pathways known to be involved in H1N1 influenza response and predicted novel targets as well. To comprehensively test our method's ability to identify potential drug targets, we used the reconstructed networks to predict the results of RNAi screens demonstrating that SDREM can make accurate predictions regarding such potential targets. SDREM correctly predicts more RNAi screen hits than top genome-wide gene prioritization algorithms, and its predictions are better suited for experimental validation

because they are placed in the context of influenza-specific pathways. We also used SDREM to predict genetic interactions that can lead to reduced viral load. Finally, we used SDREM to study the more lethal H5N1 influenza. Comparing the networks constructed for H1N1 and H5N1 allowed us to predict genes that are uniquely involved in H5N1 response, shedding light on the specific characteristics of this infection.

2 METHODS

2.1 SDREM: Signaling and Dynamic Regulatory Events Miner

SDREM (Gitter *et al.*, 2013) links the two core components of the cellular response to an external infection or treatment: signaling cascades and transcriptional regulation. To reconstruct such models, SDREM starts with the upstream proteins that detect the invading pathogen and determines signaling pathways that allow such proteins to communicate with downstream TFs. The activated TFs coordinate changes in the expression of their bound genes, which allows the cell to adapt to the new condition. SDREM integrates general TF binding data and PPI network information with condition-specific time series gene expression data to build such models.

To find TFs that are end points of the signaling pathways initiated by the upstream proteins and responsible for the observed temporal transcriptional changes SDREM iteratively combines two computational modules. For the regulatory part of the model, SDREM uses an IOHMM to analyze time series expression profiles and identify split events—places in the time series where a group of co-expressed genes starts to diverge (Ernst *et al.*, 2007). These split events are annotated with TFs that are predicted to control them allowing the method to assign temporal information to the (static) TF–gene interaction data.

Using the putative active TFs predicted for various splits, SDREM assesses the likelihood that each of these TFs is indeed responding to the infection or treatment. A TF that is not well-connected to the upstream source proteins in the signaling network is unlikely to be a major driver of the transcriptional response. To search for such connections, SDREM enumerates all possible depth-bounded paths from the sources to the TFs in the undirected PPI network. Next, SDREM orients the PPI network forcing information to flow along each edge in a single direction (Gitter *et al.*, 2011). To find the optimal orientation SDREM maximizes the following combinatorial function:

$$\sum_{p \in P} I_s(p)w(p) \quad (1)$$

where P is the set of all unique depth-bounded paths between sources and TFs, $I_s(p)$ is an indicator function that has the value 1 if path p is satisfied and $w(p)$ is the path weight. A satisfied path is a path in which all of the edges have been oriented such that the TF is reachable from the source. Our initial edge orientation analysis (Gitter *et al.*, 2011) used only edge confidence to weight paths, which in practice reduced $w(p)$ to

$$w(p) = \prod_{e \in p} w(e) \quad (2)$$

where p is a source–target path and e is an edge on the path. Edge weights $w(e)$ depend on our confidence that a specific interaction exists (Supplementary Table S1) so that more confident pathways have higher weights (Gitter *et al.*, 2011).

SDREM uses the source-TF pathways in the oriented network to quantify how well the putative active TFs are connected to the sources. This information is fed back into the temporal gene expression analysis and SDREM's two stages (IOHMM and network orientation) iterate for a fixed number of rounds. In all rounds (except for the first when pathways have not yet been learned), the IOHMM incorporates a prior that represents the likelihood each TF is activated by the signaling pathways

predicted in the previous round. In addition, the pathways can be used to identify other proteins (not TFs) that are involved in many high-weight satisfied paths and therefore important to the response.

As mentioned above, we previously used SDREM to study yeast stress response. Conceptually, SDREM is a general algorithm and can be readily applied to other organisms as well. However, when attempting to use it to model human disease response we faced several new computational challenges. First, the complexity of the human interaction networks means that SDREM's search algorithm would be too time-consuming and thus required specialized approaches to improve the runtime and accuracy. Second, while in yeast we relied only on edge and target confidence to identify high-scoring pathways, human data provides information about the nodes (proteins) as well leading to a new objective function for finding well-connected TFs. Finally, identifying key signaling proteins by determining their involvement in high-scoring pathways, as was done for yeast, may be less important for the human network. Instead, to identify potential drug targets, we would like to develop methods that can automatically determine the effects of single or double knockdowns on the ability of human cells to respond to infection or treatment. Below we discuss the new computational methods we developed to address these issues.

2.2 Algorithm parallelization and source–target pathway approximations

We explored several avenues for speeding up the computation required to reconstruct the human signaling networks. Although it is challenging to parallelize the orientation step (because local changes in edge direction may affect other edges in different parts of the network) other aspects of the algorithm could be parallelized including the randomization tests to assess TF connectivity scores and the initial depth-first search.

However, these and other precomputing steps (Supplementary Information) do not reduce the time it takes to orient the network. The Maximum Edge Orientation problem is NP-hard (Gitter *et al.*, 2011) and we already solve it using a heuristic approach, suggesting that a further approximation may not negatively impact our results. We thus modified the parallel path enumeration algorithm to only store the top m paths from any source to any TF in our dataset, ranking the paths by path weight. Considering only the top m paths also enables us to include early termination in the depth-first search's branch traversal, further reducing runtime. Evaluating the objective function requires summing the weights of all satisfied paths, and for every potential edge flip that is considered during greedy local search we must determine which paths are still satisfied. We now approximate the calculation of the orientation objective function by only examining these top m undirected paths. In test cases with millions of potential paths, the correlation between the node scores obtained using the exact objective and those obtained with the approximated objective when only using $m = 100\,000$ was >0.999 (Supplementary Figs S1 and S2). Therefore, we set $m = 100\,000$ for our H1N1 and H5N1 analysis.

2.3 Incorporating RNAi screens

When modeling human response, we can sometimes use additional sources of information regarding the involvement of a specific protein. Although in the original SDREM formulation (Gitter *et al.*, 2013) all proteins were assumed to have the same likelihood of participating in the response pathways, information such as knockdown phenotypic effects or GWAS data can increase our prior belief that a protein participates in one of the response pathways. To capture such information, we can modify the target function discussed above (Equation 2) to incorporate node priors as follows:

$$w(p) = w(t) \prod_{v \in p} w(v) \prod_{e \in p} w(e) \quad (3)$$

where p is a source–target path, t is the target on that path, v is a vertex on the path, e is an edge on the path and the function $w(*)$ is the edge confidence or node prior. Equation 3 attempts to find paths that contain proteins that are likely involved in the response based on the screen data as well as highly reliable protein interactions. Because the optimization function in Equation 2 is NP-hard (Gitter *et al.*, 2011), the target function in Equation 3 is also NP-hard and the heuristic we used to solve the original (edge only) function can be applied to this formulation as well.

Using prior information regarding protein involvement also helps us better resolve edge orientation in the signaling networks. Owing to the larger PPI and regulatory networks for human data, there are many more possible connections from sources to targets and disagreements about the orientation of individual PPI in human models. A gene that is associated with a phenotypic change in an RNAi screen leads to higher weights for the (directed) pathways in which it is contained; thus, these are preferred during the network orientation. Due to pathway redundancy, the converse is not necessarily true. Genes that are negative screen hits may still be highly relevant to signaling pathways.

Although related approaches use the screen hits directly as sources in the network (Yeger-Lotem *et al.*, 2009), we place less trust in the RNAi data. Independent RNAi screens can exhibit low overlap (Bushman *et al.*, 2009; Stertz and Shaw, 2011) in part owing to the impact of differences in methodology or cell population context (Snijder *et al.*, 2012). Supplementary Table S2 demonstrates this disagreement for RNAi screens for H1N1 influenza infection. No genes are hits in all five screens, and only a single gene is detected in four of the five screens [note that two of the screens (Bortz *et al.*, 2011; Shapira *et al.*, 2009) are targeted, not genome-wide].

Thus, to derive a prior based on screen results, we convert the RNAi data into vertex weights as follows:

$$w(v) = \begin{cases} 1 - (1 - c)^n, & \text{if } n > 0 \\ 0.5, & \text{otherwise} \end{cases} \quad (4)$$

where $w(v)$ is the weight assigned to a vertex (gene), c is the confidence in the screen in the range $[0, 1]$ and n is the number of screens that report v as a hit. We set $c = 0.75$ in all analyses here but could incorporate biological knowledge to set different confidence levels for different screens. These node priors are used directly in the formula for path weights [$w(v)$ in Equation 3] so that paths containing many screen hits have higher weights.

2.4 Predicting RNAi effects

To predict RNAi screen hits for new conditions, we developed methods to estimate the *in silico* effects of removing a protein from the signaling network component of an SDREM model. Instead of directly linking the sources and differentially expressed genes, we compute how the connectivity to the TFs is affected when a node is removed. This allows us to leverage the fact that each key TF affects many genes (often hundreds) so blocking access to such TFs significantly impacts the cell's ability to mount an effective response. We devised several scoring metrics to quantify the effect of node deletion on the targets. These metrics vary along three lines: (i) *All* versus *Top* denotes whether all satisfied paths or only the top-ranked paths are used to calculate change in connectivity. (ii) *Source–target Pairs* versus *Targets* determines whether a target's connectivity is evaluated separately for every source (i.e. each source activates a target differently) or if a target is considered to be disconnected only when it is no longer reachable from any source (i.e. any source can activate the target). (iii) *Weighted* versus *Unweighted* specifies if connectivity is treated in a binary (connected or disconnected) or continuous (fraction of connectivity remaining) fashion. The score for the

weighted variant is

$$score_w(A) = \frac{\sum_{t \in T} \frac{\sum_{p \in P(t)} w(p) I(A \notin N(p))}{\sum_{p \in P(t)} w(p)}}{|T|} \quad (5)$$

where A is the deleted node, T is the set of all targets, $P(t)$ is the set of paths to the target t to be considered (depending on the choice of All versus Top and Pairs versus Targets), $w(p)$ is the path weight, $I(*)$ is an indicator function and $N(p)$ is the set of nodes on the path. The Pairs option replaces the summation over targets by a double summation over sources and targets and updates the denominator accordingly. Intuitively, this score is the fraction of path weight that exists along paths that can still activate t after A is deleted averaged across all targets. The unweighted score, which reflects the fraction of targets that are still reachable after removing node A , is

$$score_u(A) = \frac{\sum_{t \in T} \left(1 - \prod_{p \in P(t)} I(A \in N(p)) \right)}{|T|} \quad (6)$$

2.5 Predicting genetic interactions

Although screen hits for individual genes are often noisy, higher-order knockdowns (of two or more genes) may prove to be more robust because they can target several pathways simultaneously. However, experimentally testing all possible combinations, even for pairwise interactions, is not feasible given the large number of human genes. Our method provides a viable alternative: rather than performing an experimental screen, perform *in silico* knockdowns of all pairwise combinations, score each pair and then only experimentally test the top-ranking pairs.

To derive a score for pairwise knockdowns, we mimicked experimental studies of genetic interactions in yeast (Collins *et al.*, 2010; Jonikas *et al.*, 2009; Tong *et al.*, 2004). Initial work classified gene pairs as either interacting or not with less emphasis on the strength of the interaction (Tong *et al.*, 2004). More recent work has focused on quantifying genetic interactions on a continuous scale.

$$\epsilon_{AB} = P_{AB}^{\text{observed}} - P_{AB}^{\text{expected}} \quad (7)$$

where ϵ_{AB} is the interaction between genes A and B and P_{AB} is the phenotype when both A and B are deleted (Collins *et al.*, 2010). Typically the expected phenotype is defined as the product of the phenotypes observed in the individual single deletions

$$\epsilon_{AB} = P_{AB}^{\text{observed}} - P_A^{\text{observed}} P_B^{\text{observed}} \quad (8)$$

Using this definition, negative interactions, the type we are primarily interested in because they indicate the pair blocks disease-related phenotypes, occur when the double knockdown has a stronger effect than expected because stronger effects correspond to lower values of P_{AB}^{observed} .

In yeast experimental work, colony size is a typical phenotype (Collins *et al.*, 2010; Tong *et al.*, 2004) because it approximates growth rate, but other possibilities exist (Jonikas *et al.*, 2009). Single human knockdown screens often use viral load as the phenotype (König *et al.*, 2010). In our simulations, the score defined in Equation 5 is the *in silico* phenotype P_A^{observed} . Similar to colony size and viral load, in our scoring function, more significant deletions result in lower values, and it is meaningful to take the product of the scores from two individual deletions. To calculate P_{AB}^{observed} , we generalized Equation 5 to consider the simultaneous removal of two nodes from the signaling network.

$$P_{AB}^{\text{observed}} = score_w(A, B) = \frac{\sum_{t \in T} \frac{\sum_{p \in P(t)} w(p) I(A \notin N(p)) I(B \notin N(p))}{\sum_{p \in P(t)} w(p)}}{|T|} \quad (9)$$

Equation 9 represents the average fraction of path weight that remains after removing paths that contain node A or node B .

2.6 Data

The interaction network contained PPI from BioGRID (Stark *et al.*, 2006), post-translational modifications and PPI from the Human Protein Reference Database (Mishra *et al.*, 2006) and TF-gene binding predictions (Ernst *et al.*, 2010) that were processed as described in (Schulz *et al.*, 2012) (Supplementary Information). The H1N1 data consisted of temporal gene expression (10 time points) (Shapira *et al.*, 2009), five RNAi screens (Bortz *et al.*, 2011; Brass *et al.*, 2009; Karlas *et al.*, 2010; König *et al.*, 2010; Shapira *et al.*, 2009) and 204 source proteins that detect influenza infection or directly interact with influenza proteins (Supplementary Table S3). Similarly, for H5N1, we collected temporal expression data (six time points) (Li *et al.*, 2011), one small-scale RNAi screen (Bortz *et al.*, 2011) and 41 sources (Supplementary Table S4). Source proteins were compiled from the VirHostNet database (Navratil *et al.*, 2009), large-scale host-pathogen PPI studies (Shapira *et al.*, 2009; Tafforeau *et al.*, 2011) and the literature (Supplementary Information).

We set all of SDREM's parameters to their default values (Gitter *et al.*, 2013) except the number of top-ranked paths used to score nodes and targets, which we set to 1000.

2.7 Gene prioritization algorithms

To benchmark SDREM's ability to predict H1N1 RNAi screen effects, we tested two genome-wide gene prioritization algorithms, Endeavour (Aerts *et al.*, 2006; Tranchevent *et al.*, 2008) and Pinta (Nitsch *et al.*, 2009, 2011). Endeavour ranks genes using many lines of evidence such as functional annotations, gene expression, interaction networks, text mining and sequence similarity and combines the individual rankings to create a global prioritization. Pinta is designed for diseases where a set of seed genes is not available but there are differential gene expression data between healthy and diseased individuals. It ranks genes based on the expression levels of their neighbors in an interaction network. The human proteins that interact with H1N1 proteins and differentially expressed genes were provided as input to both Endeavour and Pinta. We used EDGE (Leek *et al.*, 2006) to identify significantly differentially expressed genes by comparing the mock treatment and viral treatment time courses and using 500 null iterations. Endeavour does not use weights on the input genes so we provided the sources and all genes that had the most significant q-value from EDGE (1.64 E-6). We performed a genome-wide ranking using all lines of evidence. For Pinta the weight of a gene was $-\log_{10}(q\text{-value})$, and all genes were used as input. Sources were given the same weight as the most significantly differentially expressed genes. All default settings were used.

3 RESULTS

Pathogens infecting cells provide a clear set of sources that initiate the subsequent signaling and transcriptional response. In particular, many viruses encode only a small number of proteins allowing us to generate specific models for the host response that is triggered by host proteins that detect viral RNA or interact with viral proteins.

To test SDREM's application to such responses, we focus on influenza A viruses because of the rich datasets available and their importance to global health. The 2009 swine-origin H1N1 virus outbreak received great public attention and was declared a global pandemic in June of that year (Zhang *et al.*, 2010). More recently, research concerning mutations of avian H5N1 influenza that could allow it to be transmitted among mammals via aerosols has sparked immense controversy (Berns *et al.*, 2012),

highlighting the threat influenza A viruses pose and the need to better understand their interaction with the human immune system.

3.1 H1N1 influenza model

The host response to H1N1 influenza is the best profiled in terms of protein interactions, functional screens and transcriptional effects. We leveraged these data, in particular the temporal gene expression, to construct a dynamic model of the human immune response to H1N1 infection. SDREM identified the TFs that control this immune response and the signaling pathways that activate them—36 internal proteins that connect 33 target TFs to upstream nodes in the signaling network (Fig 1 and Supplementary Table S3). These include many proteins known to be involved in immune response such as STAT1 (Shuai and Liu, 2003), seven IRF family members (Honda and Taniguchi, 2006), three NFkB variants (Baeuerle and Henkel, 1994), RELA (Ouaaz *et al.*, 1999) and XBP1 (Martinon *et al.*, 2010). Interestingly several cancer-related proteins such as AR, BRCA1, MYC, SMAD3, SMAD7 and TP53 appear as well.

We used DAVID (Huang *et al.*, 2009) to compute the Gene Ontology (GO) (Ashburner *et al.*, 2000) biological process enrichment of the proteins predicted by SDREM, leaving out the sources because they are already known to be relevant to H1N1 and would bias the results. The most significant GO terms are dominated by processes related to transcriptional regulation due to the prevalence of target TFs in our predictions. However, beyond these are many highly relevant enriched terms including ‘immune system development’ [Benjamini-Hochberg corrected P -value 2.18 E-5 (Benjamini and Hochberg, 1995)], ‘response to virus’ ($P=0.0169$), ‘virus–host interaction’ ($P=0.0289$) and ‘immune response’ ($P=0.0431$).

3.2 Predicting RNAi screen hits

SDREM can rank human genes for their involvement in mediating host response to viral infection. Even for viruses for which genome-wide screens are available SDREM’s models are useful since they provide mechanistic explanations for genes’ involvement in the response. In addition, the large disagreement among genome-wide screens (Supplementary Table S2) makes such models important in order to distinguish real hits from noise. Finally, several pathogens, including H5N1, are challenging to work with because they require a biosafety level 3 lab, making it difficult to generate genome-wide RNAi screens. Using SDREM to predict RNAi screen hits allows researchers to prioritize candidate H5N1 targets, which could later be validated experimentally.

To determine which of the scoring metrics described in Section 2.4 is most predictive of RNAi screen hits, we ran SDREM on the H1N1 data holding out the RNAi screen data from SDREM’s input. We used the metrics to rank all 252 non-target proteins in the model using the number of high-confidence paths that use the node and the network degree to break ties in the ranking. Given the rankings for each metric, we computed the percentage of correct hits within the top X predictions (where X ranges from 10 to 100).

As seen in Table 1, SDREM performed exceptionally well on this task. To illustrate this, consider the number of correct screen hits in the top 50 ranked genes. Roughly 20 of these (depending

on the actual metric used), or 40%, are known hits. Because all the 14 334 proteins in the interaction network are included in SDREM’s search, the *expected* number of known hits in a randomly selected set of 50 genes is 3.1. Thus, using SDREM we obtain a 6.5-fold enrichment in the number of correct hits, indicating that such a method can be effectively used to design specific experiments for other viruses as well. Similar enrichments are seen for other values of X .

The best-performing metric uses only the top paths, allows targets to be activated by any source and uses the weighted score. Intuitively, including the lower confidence paths can hurt predictive performance because these source–target connections may contain false-positive PPI and not actually affect connectivity when they are broken. Our results also suggest that a target can function if it remains connected to any source. Due to the many parallel paths in the signaling network, it is uncommon for a single node removal to completely disconnect a target. The weighted variant distinguishes between deletions that do not impact a target at all and those that remove many paths to the target but do not fully disconnect it. This metric’s predictions significantly overlap with the known screen hits at all thresholds, which is also the case for most of the other metrics. Because the overlaps are significant even for the worst-performing metric when 20 or more genes are predicted, we conclude that the SDREM model itself is a powerful filter for predicting screen hits.

3.3 Comparison with gene prioritization algorithms

Having demonstrated that SDREM produces highly accurate rankings of putative RNAi screen hits, we examined whether existing tools could achieve similar performance. Although other algorithms for connecting sources and transcriptional effects provide candidate pathways (Huang and Fraenkel, 2009; Kim *et al.*, 2011; Ourfali *et al.*, 2007; Schaefer *et al.*, 2013; Tuncbag *et al.*, 2013; Yeang *et al.*, 2004; Yeager-Lotem *et al.*, 2009), they do not rank the pathway members or quantify the effect of knocking them down so it is impossible to compare our knockdown prediction results with these methods. Instead, we assess Endeavour (Aerts *et al.*, 2006; Tranchevent *et al.*, 2008) and Pinta (Nitsch *et al.*, 2009, 2011), the only two gene prioritization web servers in a recent benchmark (Börnigen *et al.*, 2012) that can prioritize the entire genome and take seed genes instead of disease keywords as input.

We ran Endeavour and Pinta using the differentially expressed genes identified by EDGE (Leek *et al.*, 2006) and the H1N1 sources as input and evaluated the number of correct screen hit predictions at the same thresholds used previously (Table 2). Because the web servers do not provide their full ranked list of genes, we cannot calculate AUC. SDREM outperforms the gene prioritization tools at all thresholds except for the top 20 threshold where Pinta makes the same number of correct predictions and Endeavour makes two more. SDREM’s advantage is greatest when assessing the top 100 predictions, at which point SDREM correctly predicts nearly twice as many RNAi screen hits as Pinta and 24% more than Endeavour. When using alternative methods to select the input genes for Endeavour and Pinta, SDREM’s strengths are even more pronounced (Supplementary Information, Supplementary Table S5). These results are especially impressive considering Endeavour uses

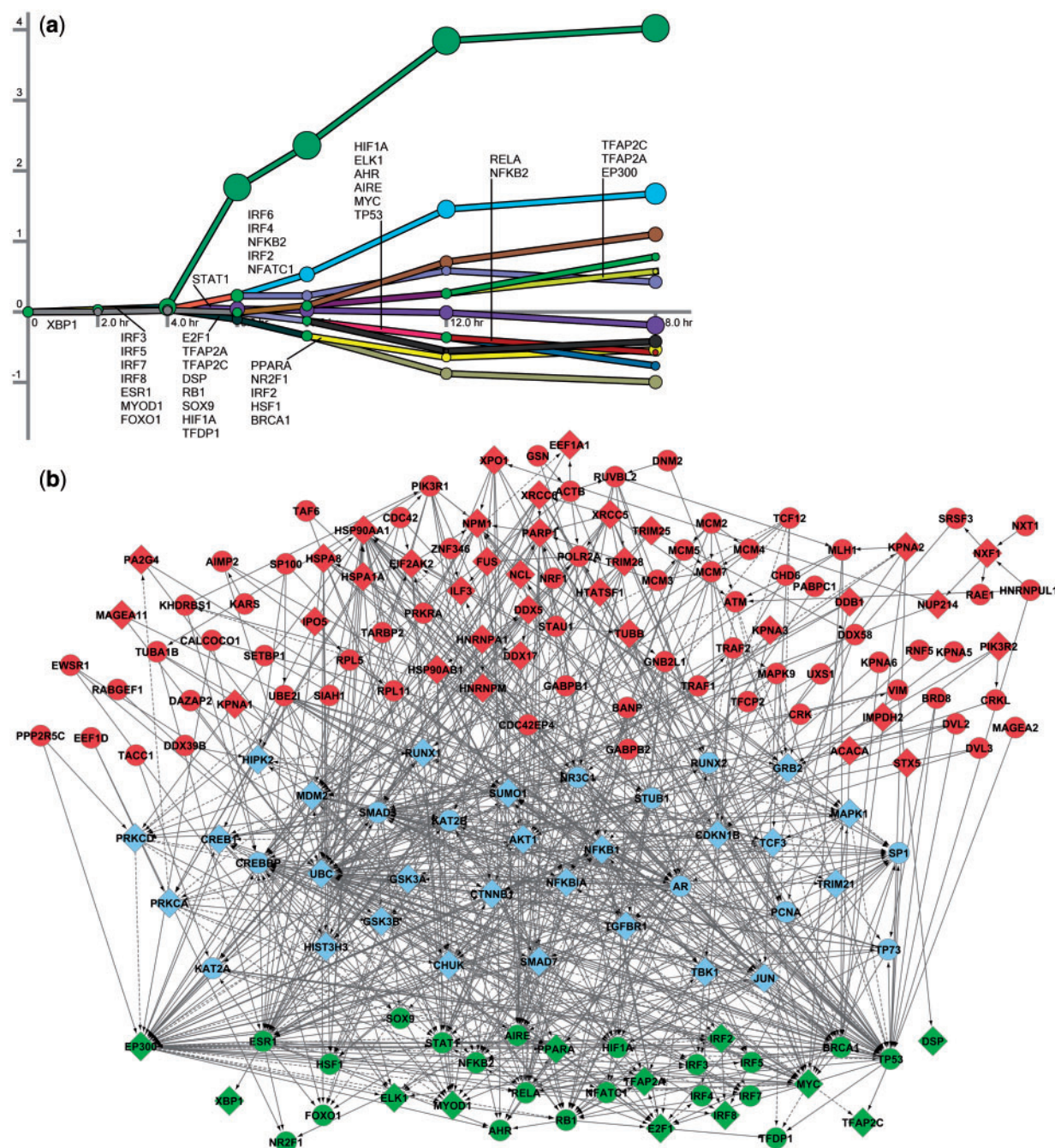


Fig. 1. The SDREM H1N1 response model. (a) The regulatory paths summarize the temporal patterns of the differentially expressed genes. The x-axis is time and the y-axis is \log_2 fold change. Split events, green nodes where a regulatory path branches, are annotated with the TFs that are predicted to activate or repress the genes at that time point. These annotations are placed on the path immediately after the split to indicate whether the TF controls the upper or lower path out of the split. (b) The signaling paths from sources (red) through internal nodes (blue) to the active TFs (green). Sources directly interact with viral proteins or detect viral presence. The active TFs are the same TFs shown on the regulatory paths. Diamonds are RNAi screen hits. Solid edges are PPIs whose orientation has been inferred by SDREM. Dashed edges are post-translational modifications and TF-gene binding interactions, which already have a known orientation

additional input that SDREM does not including prior literature via text mining. These results stress the importance of building mechanistic models that explain why differentially expressed genes are affected by infection instead of directly finding genes that are similar to the differentially expressed genes.

3.4 Predicting host proteins affecting H5N1 viral load

Having successfully used SDREM to predict H1N1 screen hits we turned to the independent H5N1 data. Using the scores from the best metric, we ranked all proteins in the H5N1 SDREM

Table 1. The scoring metrics that were used to rank H1N1 screen hits

Paths used	Connectivity	Score	AUC	Hits in top 10	Hits in top 20	Hits in top 50	Hits in top 100
Top	Targets	Weighted	0.722	6 (1.97 E-5)	8 (3.44 E-5)	18 (3.24 E-9)	42 (9.42 E-23)
Top	Pairs	Weighted	0.717	3 (2.87 E-2)	7 (2.88 E-4)	20 (4.59 E-11)	40 (8.68 E-21)
Top	Pairs	Unweighted	0.716	3 (2.87 E-2)	8 (3.44 E-5)	20 (4.59 E-11)	37 (5.59 E-18)
All	Targets	Unweighted	0.711	2 (0.153)	5 (1.09 E-2)	20 (4.59 E-11)	39 (7.82 E-20)
All	Targets	Weighted	0.706	3 (2.87 E-2)	6 (1.97 E-3)	18 (3.24 E-9)	39 (7.82 E-20)
Top	Targets	Unweighted	0.704	2 (0.153)	6 (1.97 E-3)	19 (4.02 E-10)	39 (7.82 E-20)
All	Pairs	Weighted	0.702	3 (2.87 E-2)	6 (1.97 E-3)	18 (3.24 E-9)	36 (4.43 E-17)
All	Pairs	Unweighted	0.676	2 (0.153)	6 (1.97 E-3)	18 (3.24 E-9)	36 (4.43 E-17)

Note: The metrics are sorted by area under the curve (AUC). The number of known screen hits recovered at various thresholds is shown with the significance (in parentheses) calculated using Fisher's exact test.

Table 2. Comparison of SDREM, Endeavour and Pinta gene rankings

Algorithm	Settings	Hits in top 10	Hits in top 20	Hits in top 50	Hits in top 100
SDREM	Top, targets, weighted	6	8	18	42
Endeavour	All evidence	5	10	17	34
Pinta	Default	5	8	12	22

model (Supplementary Table S4) and compared the ranks with those we obtained using the same scoring metric on the H1N1 data. We also examined the degree of the top-ranked predictions because we expected that high-scoring nodes would have high degree because such nodes are likely to affect a large number of targets when deleted. Table 3 lists the top 25 H5N1 predictions. Recall from Equation 5 that the scores denote the fraction of target connectivity remaining after the *in silico* deletion so lower scores translate to a stronger effect. Six of the H5N1 predictions—STAT3, CASP8, HSF1, ERBB3, NRIP1 and PRMT1—are particularly interesting because they are neither known H1N1 RNAi screen hits nor in the top 100 H1N1 predicted hits. Two of these, CASP8 and ERBB3, have been reported to directly interact with H5N1 viral proteins but not H1N1 proteins even though the H1N1-human PPI have much greater coverage, suggesting that they may indeed play distinct roles. Several top-ranked H5N1 proteins are sources (directly interacting with the viral proteins) or high-degree nodes. In contrast, NRIP1 and PRMT1 are neither sources nor of high degree. Their inclusion in the top predictions is noteworthy because the paths through these nodes affect targets to a greater extent than expected. NRIP1, also known as RIP140, is involved in inflammatory response as a coactivator for NF κ B (Zschiedrich *et al.*, 2008). PRMT1 methylation of STAT1 is one of the ways in which STAT1 is regulated in the immune system (Shuai and Liu, 2003).

3.5 Screen hits improve the SDREM H1N1 model

SDREM can predict screen hits, but when RNAi data are already available they do indeed lead to more accurate SDREM models. Here we refer to the model reconstructed using all data as the 'original' model and the model obtained after omitting the

RNAi data as the 'no-RNAi' model. Of the 69 signaling proteins and TFs predicted in the original SDREM H1N1 model, 38 (55%) are screen hits versus only 20 of 61 predictions (33%) in the no-RNAi model. The original and no-RNAi models overlap substantially. However, some important immune responders, including IRF2 and NF κ B1, were only detected in the original SDREM model.

The original SDREM model better corresponds to related annotated pathways as well. We used DAVID to examine the enrichment of Biocarta pathways (Nishimura, 2001) for SDREM model members. The immune-related pathways 'Human Cytomegalovirus and Map Kinase Pathways', 'The information-processing pathway at the Interferon-beta enhancer', 'T Cell Receptor Signaling Pathway' and 'Toll-Like Receptor Pathway' are enriched in the original model (all with corrected $P < 0.05$) but not the no-RNAi SDREM model.

3.6 Predicting genetic interactions

Genetic interactions are functional interactions between pairs of genes where simultaneous double deletion has a smaller or greater than expected effect. Humans have tens of thousands of genes, making it impossible to comprehensively screen for all possible genetic interactions in a condition of interest as is done to identify the phenotypic effects of single gene loss. In contrast, the *in silico* analysis performed by SDREM can be extended to pairwise or higher-order analysis. Given its performance on the single knockdown prediction task, the ranked list of pairs identified by SDREM can serve as a starting point for follow-up validation experiments. Considering the rather disappointing agreement between screen hits performed by different labs, such higher-order analysis may be required to accurately and robustly identify targets that can impact disease progression and viral load.

Based on our results for predicting H1N1 screen hits (Table 1), we again consider only the top-ranked paths and allow targets to be activated by any single source. We used our method to predict genetic interactions that affect H1N1 (Table 4) and H5N1 (Supplementary Table S6) infection. Condition-specific experimental validation is necessary to confirm that these pairs form genetic interactions that impact viral infection, but annotated pathways indicate that some of these pairs do act in parallel in other settings. For example, RB1 and TP53 are on parallel paths

Table 3. The top-ranked H5N1 RNAi screen hit predictions alongside H1N1 RNAi rankings and the number of screens reporting known hits

Gene	H1N1 source	H5N1 source	Degree	H1N1 RNAi	H5N1 RNAi	H5N1 score	H1N1 rank	H5N1 rank
HSPA8	Y	Y	95	1	1	0.765	78	1
PA2G4	Y	Y	26	1	1	0.815	66	2
AR	N	N	452	0	0	0.836	12	3
ILF3	Y	Y	39	1	1	0.901	75	4
ESR1	N	N	502	0	0	0.908	11	5
KPNA2	Y	Y	50	1	1	0.915	93	6
TP53	N	N	655	0	0	0.918	2	7
STAT3	N	N	419	0	0	0.924	151	8
CREBBP	N	N	265	0	0	0.928	53	9
SP1	N	N	365	0	0	0.931	92	10
RB1	N	N	257	0	0	0.934	5	11
GNB2L1	Y	Y	68	0	0	0.937	69	12
CASP8	N	Y	104	0	0	0.940	262	13
UBC	N	N	485	1	0	0.948	4	14
EIF2AK2	Y	Y	40	1	0	0.948	7	15
HSF1	N	N	217	0	0	0.950	N/A	16
EP300	N	N	377	1	0	0.951	3	17
BRCA1	N	N	301	0	0	0.954	49	18
NUP98	N	Y	36	2	0	0.955	N/A	19
ERBB3	N	Y	37	0	0	0.963	N/A	20
NRIP1	N	N	48	0	0	0.964	N/A	21
STAT1	N	N	642	0	0	0.964	22	22
PRMT1	N	N	70	0	0	0.964	147	23
KPNA1	Y	Y	26	1	1	0.967	216	24
HSP90AA1	Y	N	144	2	1	0.968	9	25

Note: N/A indicates that the gene was not included in the SDREM H1N1 model.

Table 4. The top 10 predicted H1N1 genetic interactions

Gene A	Gene B	ϵ_{AB}	P_{AB}^{ob}	P_{AB}^{ex}	P_A^{ob}	P_B^{ob}
EP300	TP53	−0.0077	0.8152	0.8229	0.9158	0.8986
TRAF2	UBE2I	−0.0070	0.8275	0.8345	0.9348	0.8927
UBC	UBE2I	−0.0070	0.8256	0.8326	0.9327	0.8927
RB1	TP53	−0.0068	0.8316	0.8384	0.9330	0.8986
TP53	TRAF2	−0.0066	0.8333	0.8400	0.8986	0.9348
RB1	UBE2I	−0.0057	0.8272	0.8329	0.9330	0.8927
EP300	UBC	−0.0057	0.8485	0.8541	0.9158	0.9327
EP300	TRAF2	−0.0055	0.8506	0.8561	0.9158	0.9348
EIF2AK2	UBE2I	−0.0053	0.8432	0.8485	0.9505	0.8927
NPM1	UBE2I	−0.0052	0.8442	0.8494	0.9515	0.8927

in the Biocarta pathway ‘Tumor Suppressor Arf Inhibits Ribosomal Biogenesis’ as are TP53 and TRAF2 in KEGG’s (Kanehisa and Goto, 2000) ‘Pathways in cancer’. Other predicted pairs such as ILF3 and PA2G4 (fifth on the H5N1 list) are especially interesting because the two proteins are not high-degree nodes in the PPI network.

4 DISCUSSION

SDREM is unique in that it combines time series and static data to model both signaling and dynamic regulatory networks. This

allows it to infer, more confidently, the TFs that are at the end points of signaling cascades and control human disease response. To manage the complexity of human interaction networks, we extended SDREM to leverage gene priors from RNAi screens in its objective function and incorporated several algorithmic improvements. As we have shown for influenza infection, by reconstructing disease response networks we can accurately identify key signaling pathways and nodes (proteins). SDREM can be applied successfully in many different settings even when the PPI data, node priors or gene expression time points are sparse (Supplementary Information, Supplementary Tables S7 and S8).

Given the predicted directed signaling pathways, we can estimate the phenotypic effects of knocking down genes by assessing how strongly the downstream TFs are affected. These techniques can be used to infer putative drug targets for hard-to-study conditions (such as H5N1 infection) and for combination of targets, which can lead to more robust treatments. Successfully predicting the effects of pairwise and higher-order gene knockdowns can guide targeted experimental validation, a major contribution because exhaustive pairwise RNAi screening is currently infeasible and signaling pathway redundancy can limit the effectiveness of drugs that target individual genes (Logue and Morrison, 2012). Owing to this pathway redundancy and false negatives in the existing RNAi screens, the precision we calculate for SDREM's RNAi effect predictions (Table 1) is conservative. Many proteins predicted by SDREM (for example STAT1, RELA, NFKB2 and several IRF TFs) are not screen hits but are known to be important to the immune response.

An important advantage of SDREM over previous RNAi screen and host-pathogen PPI studies is the ability to infer new pathways from general interaction data. Previous studies [for example, (Shapira *et al.*, 2009)] and pathway-based GWAS (Wang *et al.*, 2010) often rely on known curated pathways that are incomplete and not always relevant to the disease studied. In contrast, using condition-specific time series data and sources, SDREM can predict which candidate proteins are involved in the signaling pathways and which are not.

Unlike the gene scores from Endeavour, Pinta and other gene prioritization algorithms that are only defined for single genes, SDREM predicts functional disease genes by modeling the impact of removing a gene from the disease-specific signaling pathways. This simulated phenotype can be naturally extended to pairs of genes, allowing SDREM to predict genetic interactions analogously to experimental approaches. Existing genetic interaction prediction algorithms require a partial set of known genetic interactions (Bandyopadhyay *et al.*, 2012; Qi *et al.*, 2008; Wong *et al.*, 2004; Zhong and Sternberg, 2006), which prevents their application in human. Furthermore, SDREM's genetic interaction predictions are condition-specific.

Our influenza analysis focused on data from human cell lines, but some of SDREM's most exciting future applications will involve data from individual patients. We would like to use the extensions presented in this article to enable GWAS data to be used as node priors. One possible approach we intend to explore is to map SNP *P*-values to gene scores as in pathway-based GWAS (Wang *et al.*, 2010). In addition, SNPs in non-coding regions with potential regulatory functions could be used to suggest which TFs' binding is disrupted (Schaub *et al.*, 2012), providing priors for both the network orientation and the temporal expression analysis.

ACKNOWLEDGEMENT

We thank Ted Ross for his helpful discussions.

Funding: This work was supported by National Institutes of Health (1RO1 GM085022) and National Science Foundation (DBI-0965316) awards to Z.B.J. and a National Science Foundation Graduate Research Fellowship to A.G.

Conflict of Interest: none declared.

REFERENCES

- Aerts, S. *et al.* (2006) Gene prioritization through genomic data fusion. *Nat. Biotechnol.*, **24**, 537–544.
- Altshuler, D. *et al.* (2008) Genetic mapping in human disease. *Science*, **322**, 881–888.
- Ashburner, M. *et al.* (2000) Gene Ontology: tool for the unification of biology. *Nat. Genet.*, **25**, 25–29.
- Baeuerle, P.A. and Henkel, T. (1994) Function and activation of NF-kappaB in the immune system. *Ann. Rev. Immunol.*, **12**, 141–179.
- Bandyopadhyay, N. *et al.* (2012) SSLPred: predicting synthetic sickness lethality. *Pac. Symp. Biocomput.*, **2012**, 7–18.
- Bengio, Y. and Frasconi, P. (1995) An input output HMM architecture. *Adv. Neural. Inf. Process. Syst.*, **7**, 427–434.
- Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B Methodol.*, **57**, 289–300.
- Berns, K.I. *et al.* (2012) Adaptations of avian flu virus are a cause for concern. *Science*, **335**, 660–661.
- Börnigen, D. *et al.* (2012) An unbiased evaluation of gene prioritization tools. *Bioinformatics*, **28**, 3081–3088.
- Bortz, E. *et al.* (2011) Host- and strain-specific regulation of influenza virus polymerase activity by interacting cellular proteins. *mBio*, **2**, e00151–11.
- Brass, A.L. *et al.* (2008) Identification of host proteins required for HIV infection through a functional genomic screen. *Science*, **319**, 921–926.
- Brass, A.L. *et al.* (2009) The IFITM proteins mediate cellular resistance to influenza A H1N1 virus, West Nile virus, and dengue virus. *Cell*, **139**, 1243–1254.
- Bushman, F.D. *et al.* (2009) Host cell factors in HIV replication: meta-analysis of genome-wide studies. *PLoS Pathog.*, **5**, e1000437.
- Collins, S.R. *et al.* (2010) Quantitative genetic interaction mapping using the E-MAP approach. *Methods Enzymol.*, **470**, 205–231.
- Desai, K.H. *et al.* (2011) Dissecting inflammatory complications in critically injured patients by within-patient gene expression changes: a longitudinal clinical genomics study. *PLoS Med.*, **8**, e1001093.
- Ernst, J. *et al.* (2007) Reconstructing dynamic regulatory maps. *Mol. Syst. Biol.*, **3**, 74.
- Ernst, J. *et al.* (2010) Integrating multiple evidence sources to predict transcription factor binding in the human genome. *Genome Res.*, **20**, 526–536.
- Gitter, A. *et al.* (2011) Discovering pathways by orienting edges in protein interaction networks. *Nucleic Acids Res.*, **39**, e22.
- Gitter, A. *et al.* (2013) Linking the signaling cascades and dynamic regulatory networks controlling stress responses. *Genome Res.*, **23**, 365–376.
- Honda, K. and Taniguchi, T. (2006) IRFs: master regulators of signalling by toll-like receptors and cytosolic pattern-recognition receptors. *Nat. Rev. Immunol.*, **6**, 644–658.
- Huang, D.W. *et al.* (2009) Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.*, **4**, 44–57.
- Huang, S.C. and Fraenkel, E. (2009) Integration of proteomic, transcriptional, and interactome data reveals hidden signaling components. *Sci. Signal.*, **2**, ra40.
- Jonikas, M.C. *et al.* (2009) Comprehensive characterization of genes required for protein folding in the endoplasmic reticulum. *Science*, **323**, 1693–1697.
- Kanehisa, M. and Goto, S. (2000) KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*, **28**, 27–30.
- Karlas, A. *et al.* (2010) Genome-wide RNAi screen identifies human host factors crucial for influenza virus replication. *Nature*, **463**, 818–822.
- Kim, Y. *et al.* (2011) Identifying causal genes and dysregulated pathways in complex diseases. *PLoS Comput. Biol.*, **7**, e1001095.
- König, R. *et al.* (2010) Human host factors required for influenza virus replication. *Nature*, **463**, 813–817.
- Leek, J.T. *et al.* (2006) EDGE: extraction and analysis of differential gene expression. *Bioinformatics*, **22**, 507–508.
- Li, C. *et al.* (2011) Host regulatory network response to infection with highly pathogenic H5N1 avian influenza virus. *J. Virol.*, **85**, 10955–10967.
- Logue, J.S. and Morrison, D.K. (2012) Complexity in the signaling network: insights from the use of targeted inhibitors in cancer therapy. *Genes Dev.*, **26**, 641–650.
- Maher, B. (2008) Personal genomes: the case of the missing heritability. *Nature*, **456**, 18–21.
- Martinon, F. *et al.* (2010) TLR activation of the transcription factor XBP1 regulates innate immune responses in macrophages. *Nat. Immunol.*, **11**, 411–418.

- Mishra,G.R. *et al.* (2006) Human protein reference database–2006 update. *Nucleic Acids Res.*, **34** (Suppl. 1), D411–D414.
- Mohr,S. *et al.* (2010) Genomic screening with RNAi: results and challenges. *Ann. Rev. Biochem.*, **79**, 37–64.
- Moreau,Y. and Tranchevent,L.C. (2012) Computational tools for prioritizing candidate genes: boosting disease gene discovery. *Nat. Rev. Genet.*, **13**, 523–536.
- Navratil,V. *et al.* (2009) VirHostNet: a knowledge base for the management and the analysis of proteome-wide virus-host interaction networks. *Nucleic Acids Res.*, **37** (Suppl. 1), D661–D668.
- Nishimura,D. (2001) BioCarta. *Biotech Softw. Internet Rep.*, **2**, 117–120.
- Nitsch,D. *et al.* (2009) Network analysis of differential expression for the identification of disease-causing genes. *PLoS One*, **4**, e5526.
- Nitsch,D. *et al.* (2011) PINTA: a web server for network-based gene prioritization from expression data. *Nucleic Acids Res.*, **39** (Suppl. 2), W334–W338.
- Ouaaz,F. *et al.* (1999) A critical role for the RelA subunit of nuclear factor kappaB in regulation of multiple immune-response genes and in Fas-induced cell death. *J. Exp. Med.*, **189**, 999–1004.
- Ourfali,O. *et al.* (2007) SPINE: a framework for signaling-regulatory pathway inference from cause-effect experiments. *Bioinformatics*, **23**, i359–i366.
- Piro,R.M. and Di Cunto,F. (2012) Computational approaches to disease-gene prediction: rationale, classification and successes. *FEBS J.*, **279**, 678–696.
- Qi,Y. *et al.* (2008) Finding friends and enemies in an enemies-only network: a graph diffusion kernel for predicting novel genetic interactions and co-complex membership from yeast genetic interactions. *Genome Res.*, **18**, 1991–2004.
- Schaefer,M.H. *et al.* (2013) Adding protein context to the human protein-protein interaction network to reveal meaningful interactions. *PLoS Comput. Biol.*, **9**, e1002860.
- Schaub,M.A. *et al.* (2012) Linking disease associations with regulatory information in the human genome. *Genome Res.*, **22**, 1748–1759.
- Schulz,M.H. *et al.* (2012) DREM 2.0: improved reconstruction of dynamic regulatory networks from time-series expression data. *BMC Syst. Biol.*, **6**, 104.
- Shapira,S.D. *et al.* (2009) A physical and regulatory map of host-influenza interactions reveals pathways in H1N1 infection. *Cell*, **139**, 1255–1267.
- Shuai,K. and Liu,B. (2003) Regulation of JAK-STAT signalling in the immune system. *Nat. Rev. Immunol.*, **3**, 900–911.
- Snijder,B. *et al.* (2012) Single-cell analysis of population context advances RNAi screening at multiple levels. *Mol. Syst. Biol.*, **8**, 579.
- Stark,C. *et al.* (2006) BioGRID: a general repository for interaction datasets. *Nucleic Acids Res.*, **34** (Suppl. 1), D535–D539.
- Stertz,S. and Shaw,M.L. (2011) Uncovering the global host cell requirements for influenza virus replication via RNAi screening. *Microbes Infect.*, **13**, 516–525.
- Tafforeau,L. *et al.* (2011) Generation and comprehensive analysis of an influenza virus polymerase cellular interaction network. *J. Virol.*, **85**, 13010–13018.
- Tong,A.H.Y. *et al.* (2004) Global mapping of the yeast genetic interaction network. *Science*, **303**, 808–813.
- Tranchevent,L.C. *et al.* (2008) ENDEAVOUR update: a web resource for gene prioritization in multiple species. *Nucleic Acids Res.*, **36** (Suppl. 2), W377–W384.
- Tuncbag,N. *et al.* (2013) Simultaneous reconstruction of multiple signaling pathways via the prize-collecting Steiner forest problem. *J. Comput. Biol.*, **20**, 124–136.
- Wang,K. *et al.* (2010) Analysing biological pathways in genome-wide association studies. *Nat. Rev. Genet.*, **11**, 843–854.
- Wong,S.L. *et al.* (2004) Combining biological networks to predict genetic interactions. *Proc. Natl Acad. Sci. USA*, **101**, 15682–15687.
- Yeang,C. *et al.* (2004) Physical network models. *J. Comput. Biol.*, **11**, 243–262.
- Yeger-Lotem,E. *et al.* (2009) Bridging high-throughput genetic and transcriptional data reveals cellular responses to alpha-synuclein toxicity. *Nat. Genet.*, **41**, 316–323.
- Zhang,L. *et al.* (2010) Transcriptomics and proteomics in the study of H1N1 2009. *Genomics, Proteomics Bioinformatics*, **8**, 139–144.
- Zhong,W. and Sternberg,P.W. (2006) Genome-wide prediction of *C. elegans* genetic interactions. *Science*, **311**, 1481–1484.
- Zschiedrich,I. *et al.* (2008) Coactivator function of RIP140 for NFkappaB/RelA-dependent cytokine gene expression. *Blood*, **112**, 264–276.