

Systems biology

PathExt: a general framework for path-based mining of omics-integrated biological networks

Narmada Sambaturu¹, Vaidehi Pusadkar², Sridhar Hannenhalli^{3,*} and Nagasuma Chandra ^{1,2,*}

¹IISc Mathematics Initiative, Indian Institute of Science, Bangalore, Karnataka 560012, India, ²Department of Biochemistry, Indian Institute of Science, Bangalore, Karnataka 560012, India and ³Cancer Data Science Laboratory, National Cancer Institute, National Institutes of Health, Bethesda, MD 20892, USA

*To whom correspondence should be addressed.

Associate Editor: Alfonso Valencia

Received on February 12, 2020; revised on September 24, 2020; editorial decision on October 24, 2020; accepted on October 27, 2020

Abstract

Motivation: Transcriptomes are routinely used to prioritize genes underlying specific phenotypes. Current approaches largely focus on differentially expressed genes (DEGs), despite the recognition that phenotypes emerge via a network of interactions between genes and proteins, many of which may not be differentially expressed. Furthermore, many practical applications lack sufficient samples or an appropriate control to robustly identify statistically significant DEGs.

Results: We provide a computational tool—PathExt, which, in contrast to differential genes, identifies differentially active paths when a control is available, and most active paths otherwise, in an omics-integrated biological network. The sub-network comprising such paths, referred to as the TopNet, captures the most relevant genes and processes underlying the specific biological context. The TopNet forms a well-connected graph, reflecting the tight orchestration in biological systems. Two key advantages of PathExt are (i) it can extract characteristic genes and pathways even when only a single sample is available, and (ii) it can be used to study a system even in the absence of an appropriate control. We demonstrate the utility of PathExt via two diverse sets of case studies, to characterize (i) *Mycobacterium tuberculosis* response upon exposure to 18 antibacterial drugs where only one transcriptomic sample is available for each exposure; and (ii) tissue-relevant genes and processes using transcriptomic data for 39 human tissues. Overall, PathExt is a general tool for prioritizing context-relevant genes in any omics-integrated biological network for any condition(s) of interest, even with a single sample or in the absence of appropriate controls.

Availability and implementation: The source code for PathExt is available at <https://github.com/NarmadaSambaturu/PathExt>.

Contact: nchandra@iisc.ac.in or sridhar.hannenhalli@nih.gov

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Whole-genome transcriptomic data are routinely harnessed to probe genes and processes underlying specific biological contexts, including diseases (Jiang *et al.*, 2015). Extracting biological insights from such high-dimensional data remains an important challenge (Esteve-Codina, 2018). A standard approach to interpreting such data is to first identify differentially expressed genes (DEGs) and then to identify enriched functions among such genes (Esteve-Codina, 2018). However, biological phenotypes emerge from complex interactions among numerous biomolecules, resulting in a highly heterogeneous transcriptional landscape, thus adversely affecting the power to detect critical genes and pathways based on DEGs alone. Moreover, such high-coverage data encodes a vast amount of information

beyond DEGs, warranting exploration using multiple complementary approaches. Genome-wide molecular interaction networks constructed from experimentally identified physical, regulatory, signaling and metabolic interactions have shown great promise as a framework for integrating and interpreting such data (Sambarey *et al.*, 2017a,b). The identification of sub-networks in such biological networks, which encode the processes perturbed by a stimulus, or active processes in general, can lead to mechanistic insights, as well as help prioritize genes for intervention (Mitra *et al.*, 2013).

Several methods have been proposed to integrate transcriptomic data with biological networks, that identify ‘active modules’ or connected sub-networks which show changes across conditions (Mitra *et al.*, 2013). Current approaches are largely built on the work by Ideker *et al.* (2002), called jActiveModules, which formulates a sub-

network scoring scheme based on the statistical significance of differential gene expression, and then identifies high-scoring sub-networks using a simulated annealing approach. Other methods along similar ideas have been proposed, that filter sub-networks based, for example, on network motifs (Milo *et al.*, 2002), or on k -shortest paths between a set of ‘seed’ nodes sampled based on their differential expression (Cabusora *et al.*, 2005). He *et al.* (2011) study the dynamics in hepatocellular carcinoma by identifying an active sub-network for each stage of the disease by only retaining edges linking statistically significant DEGs, and then comparing the different sub-networks. Despite the availability of interaction data, these methods largely rely on network scoring schemes which prioritize DEGs (Mitra *et al.*, 2013). However, in many practical scenarios including clinical settings, lack of appropriate controls or sufficiently large number of samples preclude robust identification of statistically significant DEGs (Stretch *et al.*, 2013).

In this work, to complement the conventional differential expression-based analyses, we provide PathExt, a path-based approach to mining omics-integrated biological networks. PathExt uses a network weighting scheme that prioritizes edges/interactions rather than nodes/genes, and identifies differentially active paths when comparing conditions, or highly active paths when studying a single condition. The sub-network comprised of these differential paths, referred to as the TopNet, captures the genes and pathways characterizing the biological condition under study. Deviating from traditional approaches to active sub-network identification, PathExt does not use the selection of a connected module as a constraint. Rather, the method results in a well-connected sub-network, reflective of the interconnectedness of biological processes responding to any stimulus.

PathExt can be used to address the following biologically important questions: (i) What are the most significantly differential paths between conditions, and what are the most critical genes underlying the differentially active paths (note that the critical genes themselves may not be differentially active)?; (ii) What is the response to a given perturbation?; and (iii) What are the most active paths and processes in a condition for which there is no appropriate control?

We demonstrate the wide applicability of PathExt by applying it to two diverse sets of case studies. (i) Exposure of the pathogen *Mycobacterium tuberculosis* to 18 antibacterial drugs, where only one sample is collected for each such exposure. We find that the TopNet for each sample reveals the pathways known to be affected by the corresponding drug. (ii) Transcriptomic data for 39 human tissues. Application of PathExt reveals tissue-relevant genes and processes despite the absence of a clear control. In all applications, we find that the TopNet forms a well-connected graph (not expected by chance). Overall, PathExt is a general framework for the integration and analysis of knowledge-based biological networks and omics data, to reveal context-relevant genes and processes. This can be done even with a single sample, or in the absence of appropriate controls. We provide the open source PathExt tool at <https://github.com/NarmadaSambaturu/PathExt>.

2 Materials and methods

2.1 PathExt

We provide an overview of PathExt in Figure 1. The inputs to PathExt are (i) a directed gene network and (ii) gene-centric omics data for the conditions of interest. The omics data can represent a variety of quantities pertaining to the node, such as gene expression level, differential expression, protein, metabolite level, etc., in one or more conditions. The output of PathExt is a sub-network, that we refer to as the TopNet, consisting of the most significant differential or active paths, and is interpreted based on the application context.

PathExt can be used to interrogate any combination of knowledge-based networks and omics data. For clarity, we describe the steps for a protein-protein interaction network (PPIN) and gene expression data. The pipeline consists of the following steps (Fig. 1):

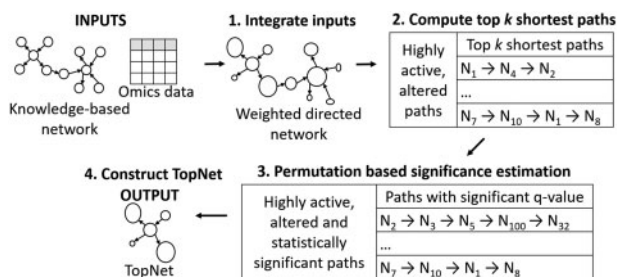


Fig. 1. PathExt overview. PathExt uses a knowledge-based directed network and omics data as inputs, and outputs a sub-network consisting of context-relevant genes and processes, referred to as the TopNet. 1. Integrate inputs: PathExt integrates the inputs by weighting the nodes and edges of the knowledge-based network as a function of the abundance of the biomolecules. Node weight is one of $SI \times FC$, FC or SI , where SI is signal intensity (e.g. gene expression) and FC is the fold change in abundance. Edge weight $_{(i,j)} = \frac{1}{\sqrt{N_i \times N_j}}$, giving an edge between highly abundant or altered molecules a low weight. Thus, a low-weight path will traverse highly active or altered molecules a low weight. 2. Compute top k shortest paths: Shortest paths are computed between all pairs of nodes in the weighted network, of which the top k paths with least weight, corresponding to highly active and altered paths, are used for further analyses. 3. Permutation based significance estimation: The rows of the input omics matrix are shuffled independently to create r random weighted networks, and estimate statistical significance of the top k shortest paths. 4. Construct TopNet: Among the top k paths, the multiple testing corrected significant paths constitute a sub-network, called the TopNet

(1) integrate inputs, (2) compute top k shortest paths, (3) permutation based significance estimation and (4) construct TopNet by retaining the edges in the significant shortest paths.

1. Integrate inputs: We integrate the inputs by computing (sample-specific or condition-specific) node and edge weights in the knowledge-based network using the omics data. In the specific scenario when comparing conditions (e.g. pre- and post-treatment), we encode the ‘response’ of the system to the change in conditions by assigning the node weight as either the fold change in gene expression ($N_i = FC$), or fold change in combination with gene expression ($N_i = SI \times FC$). Here, N_i is the weight of node i , and SI is the normalized signal intensity, or expression level, of a particular gene. Assigning $N_i = SI \times FC$ prioritizes abundant genes which may not be highly differentially expressed, and also distinguishes between two genes with the same fold change, but different abundance. For instance, given $FC(\text{gene A}) = 400/200$, and $FC(\text{gene B}) = 4/2$, genes A and B will be treated equally if $N_i = FC$, but gene A will be prioritized if $N_i = SI \times FC$. Thus, using $N_i = FC$ can be expected to highlight alterations between conditions, while $N_i = SI \times FC$ can help study highly active as well as altered processes. We have provided both options in our tool for the user to choose from depending on the application. The response to a perturbation can be studied in terms of up-regulated/activated pathways (Activated Response TopNet), obtained by computing $FC = SI_{\text{perturbed}}/SI_{\text{control}}$, or down-regulated/repressed pathways (Repressed Response TopNet), obtained using $FC = SI_{\text{control}}/SI_{\text{perturbed}}$. The Response TopNet is a union of these two TopNets, and provides a holistic view of the active, altered genes and processes. Exclusively applying the expression value as the node weight ($N_i = SI$) is useful either when no control is available, or when the emphasis is on identifying highly active processes in each state. This TopNet is referred to as the Highest Activity TopNet (HA TopNet). Even in this case, comparisons between states can be carried out after the TopNet is generated for each state.

We interpret an edge to represent a ‘reaction’ between the two nodes, and following the principles of mass action kinetics, an edge between highly abundant nodes is given Edge weight $_{(i,j)} = \frac{1}{\sqrt{N_i \times N_j}}$, where N_i and N_j are the weights of the incident nodes i, j . This choice gives highly active interactions a low edge weight.

2. Compute top k shortest paths: To achieve a biological outcome, typically a sequence of active reactions is involved, represented by a series of low weight edges in our network. In order to enumerate such low

weight paths, we use Dijkstra's algorithm (Dijkstra, 1959) to identify all-pair-shortest-paths. We then normalize the path weight for each node pair by the number of edges along the shortest path to get Normalized path weight = $\frac{\sum_{\text{edges in path}} \text{Edge weight}}{\text{Number of edges in path}}$, and retain the top k shortest paths with least weight. Here k is a user-defined, application-specific threshold.

3. Permutation based significance estimation: We assess the statistical significance of the normalized weight of each selected path empirically as follows. Given an $m \times n$ matrix of gene expression data for m genes in n samples/conditions, we randomly shuffle data in each row (gene) independently. The edges are re-weighted with the randomized gene expression data, and the weight of each path from step 2 is computed. After r such randomizations, for each path selected in step 2, r randomization-based weights are computed, based on which a z -score and p -value is estimated for each path. The p -value is finally transformed into a q -value (Benjamini and Hochberg, 1995) to account for multiple hypotheses testing. All paths with significant q -value are retained.

4. Construct TopNet: The edges in the significant paths from step 3 form a sub-network, which we refer to as the TopNet. The TopNet provides a snapshot of the active and/or significantly altered processes in the system, and can be studied to gain mechanistic insights. To further prioritize critical genes and paths in the TopNet, we apply network centrality measure—Ripple Centrality (Sambaturu et al., 2016).

In cases where a single condition is being examined, or the number of conditions is too small to generate a sufficiently large number of randomized gene expression matrices, step 3 can be skipped, and top k shortest paths can be taken to represent highly active, altered paths, albeit without the statistical filter. In such cases, Step 4 can be directly applied to these paths to generate a TopNet.

2.2 Ripple centrality

Ripple centrality (Sambaturu et al., 2016) prioritizes nodes which can reach a large fraction of the network along highly active and perturbed paths. It is measured as Ripple centrality(u) = $C(u) \times R_{\text{out}}(u)$, where $R_{\text{out}}(u) = |\text{nodes reachable from } u|$ denotes the outward reachability of node u , and $C(u) = (n-1) / \sum_{v=1}^{n-1} \sigma(u,v)$ gives the closeness centrality of node u . Here $\sigma(u,v)$ denotes the weight of the shortest path from node u to all $n-1$ other nodes in the graph.

2.3 Mycobacterium tuberculosis (M.tb) drug exposure

2.3.1 Data

Transcriptomic data for M.tb H37Rv exposed for 16 h to 2xMIC of 18 drugs was obtained from GSE71200 (Ma et al., 2015). The list of 18 drugs along with their mechanism of action and TopNet details can be found in Supplementary Table S1. A knowledge-based network composed of experimentally validated protein-protein interactions as well as regulatory interactions in M.tb was obtained from Mishra et al. (2017), consisting of 3686 genes and 34 223 edges.

2.3.2 Gold standards

INH is known to affect the mycolic acid synthesis and processing pathways in M.tb (Wishart et al., 2018). To create a gold standard for INH treatment, we searched for the term 'mycolic acid' in Mycobrowser (Kapopoulou, 2011), a database of manually curated annotations for pathogenic mycobacteria, including M.tb. This resulted in a list of 17 M.tb genes, to which we added katG and fas, the known targets of INH (Wishart et al., 2018). Similarly, gold standards were created for 5 other drugs by searching for terms related to their known mechanisms of action—'RNA polymerase' for Rif, 'mycolic acid' for ethionamide, 'protein synthesis' for capreomycin and '30s ribosomal protein' as well as '16s rRNA' for kanamycin and streptomycin (Wishart et al., 2018) (Supplementary Table S3).

2.3.3 TopNet creation

For all 18 drugs in GSE71200 (Ma et al., 2015), Activated Response TopNets were constructed using $N_i = SI_{\text{drug}} \times (SI_{\text{drug}}/SI_{\text{control}})$, while $N_i = SI_{\text{control}} \times (SI_{\text{control}}/SI_{\text{drug}})$ was used to construct the Repressed Response TopNets. Only shortest paths with 2 or more edges were considered, and 1000 randomizations of the gene expression matrix were carried out for computing statistical significance of shortest paths. The percentile and q -value thresholds were chosen such that the resulting TopNets were of similar size for all cases (Supplementary Table S1). Activated and Repressed TopNets are provided in Supplementary Files S1 and S2, respectively.

2.3.4 Functional enrichment

Functional enrichment was carried out using ClueGO v2.3.4 (Bindea et al., 2009), a plugin in the network visualization tool Cytoscape 3.2 (Shannon et al., 2003). Enrichment was against GO Biological Processes, GO Cellular Components and GO Molecular Functions, with a q -value cutoff of 0.01. Enriched pathways for all 18 drug exposure cases are provided in Supplementary File S3.

2.3.5 Significance of TopNet connectedness

Significance of TopNet connectedness was tested by comparing against comparable sub-networks induced by (i) the top DEGs, (ii) 1000 sets of randomly sampled genes and (iii) 1000 sets of randomly sampled edges. Here the number of DEGs and sampled genes (or edges) corresponds to the number of nodes (or edges) in the TopNet.

2.3.6 Comparison with existing methods

We compared PathExt with jActiveModules (Ideker et al., 2002) and the method developed in He et al. (2011), referred to here as *DEG networks*, for the 6 M.tb drug exposure cases with gold standards. Results were evaluated by comparing with these gold standards. Typically, both these methods depend on the statistical significance of DEGs to identify an active sub-network. This could not be computed for the data considered here as only 1 sample was available per condition.

We created DEG networks by considering all genes 1.5-fold up- or down-regulated as DEGs, and retaining edges linking these DEGs.

The jActiveModules plugin v3.2.1 in Cytoscape 3.8 (Shannon et al., 2003) allows users to provide inputs other than p -value, and specify whether larger or smaller values are to be considered most significant. We provided $|\log_2(FC)|$ as input, with the specification that large values be considered significant. jActiveModules was also executed with the raw fold change values.

2.4 Human tissues

2.4.1 Data

Normalized gene expression data was collected from GTEx (Carithers and Moore, 2015) (dbGaP accession number phs000424.v7.p2) for 39 human tissues, corresponding to 23 organs and 2 cell lines. The signal intensities of each tissue were summarized using the LMFfit function in R (Limma package; Ritchie et al., 2015). The antilog of the fitted value was used for further analysis as PathExt requires non-negative values. Human protein-protein interaction network (hPPIN) comprising regulatory, signaling and metabolic pathways was obtained from (Sambarey et al., 2017a). This network has 17 062 proteins (nodes) and 208 759 interactions (edges).

2.4.2 TopNet creation

Since no control was available, we constructed two types of TopNets—HA TopNets using $N_i = SI$, and z -score TopNets using $N_i = |z - \text{score}|_i$. Here, z -score for a gene i in a given tissue was computed with respect to all tissues, and statistical significance of shortest paths was computed by randomizing the $|z - \text{score}|$ matrix 1000 times. The size of the TopNet can vary across tissues and across percentile and false discovery rate thresholds. For the z -score

TopNets, we explored percentile thresholds in the range [0.001, 1.0] and q-value thresholds from the set {0.001, 0.005, 0.01, 0.05} in each tissue to adjust the TopNet size to ≈ 300 nodes. Then for the HA TopNet of each tissue, we explored the same set of percentile thresholds so as to have a comparable size between HA and z-score TopNets; the percentiles across tissues were either 0.001 or 0.002 in all cases. Thresholds for all tissues are available in [Supplementary Table S9](#). HA TopNets and z-score TopNets for all tissues are provided in [Supplementary Files S4 and S5](#), respectively.

2.4.3 Gold standards

The human protein atlas (HPA; [Uhlén et al., 2015](#)), a compiled list of Disease genes ([Feiglin et al., 2017](#)) and genes from the Disease Ontology browser of the Mouse Genome Informatics (MGI) database ([Bult et al., 2019](#)) were used to validate the results. HPA provides lists of genes whose mRNA expression is elevated in a particular tissue. The elevated expression can correspond to one of three categories: (i) ≥ 5 -fold mRNA levels in a particular tissue as compared to all other tissues, (ii) ≥ 5 -fold higher mRNA levels in a group of 2–7 tissues and (iii) ≥ 5 -fold higher mRNA levels in a particular tissue as compared to average levels in all tissues. The union of genes from the above three categories form the gold standard. HPA data was downloaded on the 26th of December, 2018. Disease genes were compiled by [Feiglin et al. \(2017\)](#) by cross-referencing data from two databases—Online Mendelian Inheritance in Man (OMIM, [Hamosh et al., 2004](#)), and the Human Phenotype Ontology (HPO, [Köhler et al., 2014](#)). OMIM is a compendium of associations between genetic variations and predominantly Mendelian disorders, while HPO provides a standardized vocabulary for working with such phenotypic abnormalities. The Disease Ontology browser of the MGI lists genes whose mutation is associated with phenotypes characteristic of human disease ([Bult et al., 2019](#)). A list of housekeeping genes obtained from [Eisenberg and Levanon \(2013\)](#), comprising of 3804 genes with constant expression level across a panel of tissues, is used as a negative control to test whether tissue TopNets are enriched in ubiquitously active genes.

2.4.4 Functional enrichment and ranking of pathways

Enrichment was carried out using the `enrichGO` function of the R package `clusterProfiler` v3.6.0 ([Yu et al., 2012](#)), using Biological Processes as the ontology, and with a Benjamini Hochberg cutoff of 0.01. For each tissue, the background for enrichment was set to be the list of genes for which both expression and interaction data were available. Pathway enrichment results for HA TopNets, z-score TopNets, their corresponding baselines, gold standards, as well as housekeeping genes, are provided for all tissues in [Supplementary File S6](#). Pathways enriched in each tissue which are enriched in at most 10% of tissues (≤ 4 tissues) are listed in [Supplementary Tables S12](#) (HA TopNet) and [S13](#) (z-score TopNet). Pathways enriched in the TopNets were ranked based on the weight of the first shortest path involving a gene from that pathway. Ties were broken based on the fold enrichment of TopNet genes in a pathway relative to expectation.

3 Results

3.1 PathExt reveals pathways related to drugs' mechanism of action in treated M.tb

In a previous study, the *Mycobacterium tuberculosis* (M.tb) strain H37Rv was exposed to different concentrations of various anti-tuberculosis drugs, and the transcriptional response was measured (GEO accession number GSE71200; [Ma et al., 2015](#)). We obtained the transcriptomic data for 2xMIC (twice the minimum inhibitory concentration) dose of 18 drugs, for bacteria surviving 16 h of drug exposure, suggesting a degree of drug resistance. Only one replicate per MIC per drug and a single untreated control sample were measured, making robust estimation of differential expression impractical. For 6 drugs where the mechanism of action is well studied ([Wishart et al., 2018](#)), we obtained gold standard sets of genes

experimentally verified to be perturbed upon drug exposure (Section 2.3.2). In all 6 cases, the Response TopNets generated by PathExt are concordant with the gold standards, and reveal genes and pathways relevant to the action of each drug ([Table 1](#)). In contrast, the genes with 1.5-fold differential expression have consistently poor overlap with gold standards ([Table 1](#)). We discuss the Isoniazid and Rifampicin exposures in detail below.

3.1.1 PathExt links INH exposure to mycolic acid synthesis and processing

The anti-bacterial drug Isoniazid (INH) inhibits the synthesis of mycolic acids, which are long fatty acids found in the cell walls of mycobacteria ([Wishart et al., 2018](#)). The Activated Response TopNet (selecting for up-regulated paths), Repressed (down-regulated paths) and merged Response TopNets (Section 2) identified by PathExt were all significantly enriched in gold standard genes related to mycolic acid synthesis and processing ([Table 1](#), [Supplementary Table S1](#)). In stark contrast, the DEGs with ≥ 1.5 -fold differential expression had poor overlap with the gold standard ([Table 1](#), [Supplementary Table S1](#)). The central genes (Section 2.2) in the Activated Response TopNet consist of genes involved in mycolic acid biosynthesis, whereas the Repressed Response TopNet has unsaturated acyl-CoA hydratases responsible for oxidizing fatty acids, and genes involved in lipid degradation as the central nodes ([Supplementary Tables S4 and S5](#)). These results unambiguously point to the up-regulation of fatty acid synthesis and down-regulation of its degradation as a resistance response to INH exposure.

A previous study ([Takayama et al., 2005](#)) consolidated experimental and computational evidence to list the 7 main processes in the mycolic acid synthesis and processing pathway, namely, the FAS-I (fatty acid synthetase-I) system, transition from the FAS-I system to the FAS-II system, the FAS-II system, cyclopropane synthases and methyltransferases, oxidation-reduction, Claisen-type condensation and mycolic acid processing. Of the 42 genes described in their work, interaction and expression data were available for 39, of which 16 were present in the INH exposed Response TopNet (3.63 fold enrichment; Fisher's p -value = $1.68e-6$), while the 1078 DEGs comprise only 14 of these genes (Fisher's p -value = 0.27). Notably, the TopNet sub-network induced by the 16 genes from the mycolic acid synthesis and processing pathway ([Fig. 2](#)) and their immediate neighbors, represent all 7 component processes. Interestingly, NADH dehydrogenases (highlighted in violet in [Fig. 2](#)) are also picked up in this sub-network. It has been hypothesized that M.tb may gain resistance to INH by regulating NADH dehydrogenase and the intracellular NADH/NAD⁺ ratio ([Miesel et al., 1998](#)). This is consistent with the fact that the bacteria under study are the ones which survived exposure to 2xMIC of INH and thus likely to have triggered their resistance processes.

Finally, as an additional control, we directly compared the Response TopNet genes with same number of top DEGs in terms of their functional enrichment ([Fig. 2](#), Section 2.3.4). The genes in the Response TopNet are enriched in the functional terms relevant to INH exposure, such as *cell periphery*, which is the part of the cell most affected by INH ([Wishart et al., 2018](#)), and stress response terms such as *oxidoreductase activity* and *oxidation-reduction process*. We also find the term *regulation of metabolic processes*, which is an expected energy conservation response. In contrast, the top 401 DEGs are enriched for the terms *quimone binding* and *symbiosis encompassing mutualism through parasitism*, which are not informative of the condition under study. Together, these results show that the Response TopNet for M.tb exposure to 2xMIC of INH is indeed characteristic of its action and reveals genes and processes that would be missed by a conventional approach relying on differential gene expression alone.

3.1.2 Rif exposure TopNet reveals the perturbation of nucleotide synthesis pathway

Rifampicin (Rif) inhibits DNA-dependent RNA polymerase activity, thus suppressing transcriptional initiation ([Wishart et al., 2018](#)).

Table 1. Evaluation of TopNets

Drug	Accession number	Gold standard		Activated Response TopNet		Repressed Response TopNet		Response TopNet		1.5 FC DEGs			
		Nodes	P-value standard	Nodes	P-value standard	Nodes	P-value standard	Nodes	P-value standard	DEGs	DEGs in base network	Gold standard	P-value
Capreomycin	GSM1829654	29	0.45	184	0.02	195	0.02	372	0.03	1796	1676	16	0.26
Ethionamide	GSM1829659	15	0.001	201	1	202	1	394	0.02	610	562	6	0.02
Isoniazid	GSM1829740	17	0.0002	195	0.0002	213	0.0002	401	2.2e-7	1078	987	8	0.07
Kanamycin	GSM1829743	30	0.02	191	0.5	199	0.5	379	0.03	790	731	4	0.23
Rifampicin	GSM1829752	22	0.0001	196	0.03	197	0.03	380	4.3e-6	1579	1462	13	0.07
Streptomycin	GSM1829755	30	0.0009	193	0.06	181	0.06	338	0.0003	308	286	1	0.29

Note: Response TopNets for M.tb exposed to six drugs whose mechanism of action is well known, are concordant with gold standards and reveal genes relevant to the action of each drug. A standard DEG analysis shows poor concordance.

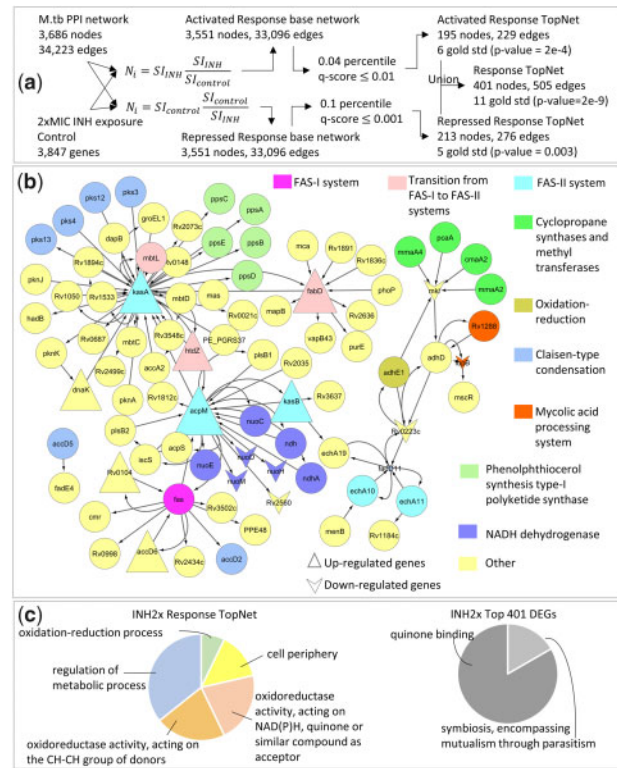


Fig. 2. Response to 2xMIC INH. (a) Gene expression data for a single sample of M.tb exposed to 2xMIC of INH for 16 h is integrated with a knowledge-based protein-protein interaction network for M.tb. (b) Sub-network of the Response TopNet formed by extracting genes from the mycolic acid synthesis and processing pathway (Takayama et al., 2005), the known target pathway of INH and their immediate interactors. Every component process of this pathway is represented in the Response TopNet by at least 1 gene. (c) GO enrichment of Response TopNet gives pathways relevant to INH exposure, such as cell-periphery and oxidation-reduction process. Enrichment of an equal number of top DEGs does not provide drug-specific insights

Once again, the Activated, Repressed and union Response TopNets are enriched in gold standard genes, whereas the DEGs are not (Table 1, Supplementary Table S1). The gene rpoB (Rv0667) is central in the Activated Response TopNet, effectively recapitulating previous reports which suggest that Rif resistance can be caused by transcriptional up-regulation of rpoB (Zhu et al., 2018). The error prone DNA repair synthesis protein DnaE2 (Rv3370c), and the genetic recombination and nucleotide excision repair protein RecA (Rv2737c) are also central in this network. Exposure to antibiotics such as Rif has been shown to result in a recA-dependent SOS response, and a corresponding increase in dnaE2 levels (McGrath et al., 2014). Also, the up-regulation of dnaE2 has been identified as a critical factor in the emergence of drug resistance both *in-vitro* and *in-vivo* (Boshoff et al., 2003). Other central genes (full list in Supplementary Table S4) include the 16S ribosomal RNA methyl-transferase Rv2372c, and the replicative DNA helicase dnaB (Rv0058). These genes reflect perturbations in the nucleotide synthesis pathway, the very pathway known to be affected upon exposure to Rif. Central genes in the Repressed Response TopNet include, among others, dnaK (Rv0350) and Rv0232, a transcriptional regulator of the tetR/acrR-family. Disruption of Rv0232 has been shown to provide a growth advantage to H37Rv *in-vitro* (DeJesus et al., 2017). We found that Rv0232 was 4.5-fold down-regulated and centrally involved in repressed paths, suggesting this as a possible resistance mechanism.

Interestingly, dnaK is central in the Repressed Response TopNet for Rif, whereas it is central in the Activated Response TopNet for INH exposure. It has been shown that dnaK is repressed by Rif (Eltringham et al., 1999), whereas cells with higher levels of dnaK are more likely to persist upon exposure to INH (Jain et al., 2016).

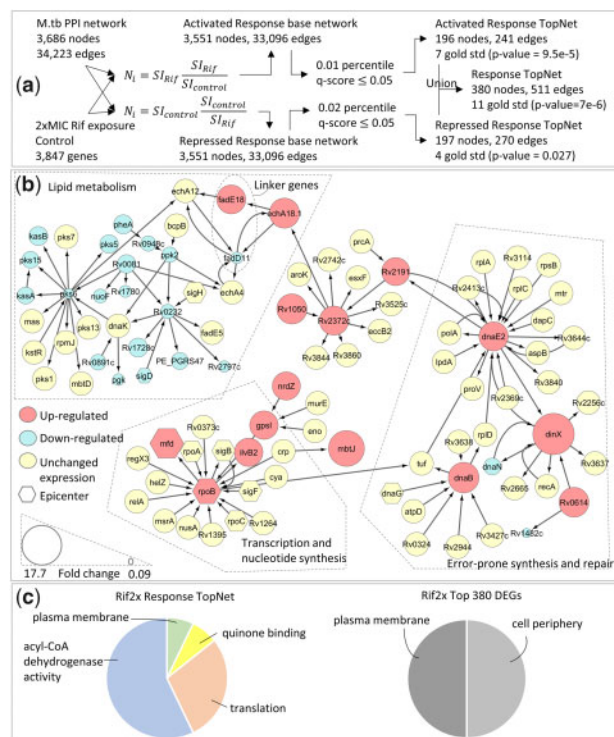


Fig. 3. Response to 2xMIC Rif. (a) Gene expression data for a single sample of M.tb exposed to 2xMIC of Rif for 16 h is integrated with a knowledge-based protein-protein interaction network for M.tb. (b) Central genes from Activated and Repressed Response TopNets along with their immediate interactors, extracted from the union Response TopNet for Rif. This module contains genes related to transcription and nucleotide synthesis, the known target pathway of Rif. Other pathways represented here are lipid metabolism and error-prone synthesis and repair, both known mechanisms of resistance to Rif (Boshoff *et al.*, 2003; Howard *et al.*, 2018). (c) GO enrichment of Response TopNet gives pathways relevant to Rif exposure, such as translation. Enrichment of an equal number of top DEGs does not provide drug-specific insights

This result underscores the biological and mechanistic relevance, as well as the condition-specificity of the TopNets generated by PathExt.

Although the exact pathway for DNA-dependent RNA polymerase activity is not known, examining the central genes from the Rif Activated and Repressed Response TopNets along with their immediate interactors provides valuable insights. These genes form two connected components, connected by two linker genes, *fadE18* (Rv1933c) and *fadD11* (Rv1550) (Fig. 3). This sub-network highlights three major error processes, namely, (i) transcription and nucleotide synthesis, (ii) error-prone synthesis and repair and (iii) lipid metabolism. Figure 3 also shows the GO-term based enrichment of the genes in the Response TopNet, and for an equal number of top DEGs. The genes in the Response TopNet are enriched for terms relevant to exposure to Rifampicin, such as *translation*, which is the process targeted by Rif, *plasma membrane* and *acyl-CoA dehydrogenase activity*, which are related to lipid metabolism. On the other hand, the 380 top DEGs are enriched for the terms *cell periphery* and *plasma membrane*, which are not specifically informative of cellular response to the drug.

As demonstrated by the INH and Rif case studies, each Response TopNet reveals drug-specific mechanisms. Drug-specificity of the TopNets is further emphasized by the fact that there is no node or edge common to all 18 Response TopNets, despite the same knowledge-based network being used as input in all cases.

The Response TopNet is a connected graph with > 50% nodes in the largest component in each of the 18 drug exposures. This connectedness, reflective of biological pathways, is shown to be non-random (Section 2.3.5), and not captured by the sub-networks induced by top DEGs (Supplementary Table S2). This suggests that

our Response TopNet captures crosstalk between the dysregulated paths, which simple differential gene expression analysis may not.

Taken together, these results show that PathExt captures drug-specific responsive genes and processes, even when only a single sample was available per condition.

3.1.3 PathExt outperforms existing DEG-based methods

Table 2 compares PathExt with jActiveModules and DEG networks, highlighting some of their distinguishing characteristics and comparing the sub-networks identified for each of the 6 M.tb drug exposure cases with gold standards. The Response TopNet provided by PathExt is always enriched with gold standard genes, while the modules identified by jActiveModules with input $|\log_2(FC)|$ and DEG networks are enriched in only 2 and 3 cases, respectively. When provided with raw FC values as input, the networks generated by jActiveModules are enriched with gold standards in 3 cases (Supplementary Table S1). This shows that PathExt outperforms traditional DEG-based methods in these cases, and highlights its value particularly when there is a paucity in the number of samples.

3.2 Human tissue TopNets reveal tissue-related genes and processes

In a second set of case studies, we applied PathExt to identify tissue-related pathways using gene expression data for 39 human tissues in GTEx (Carithers and Moore, 2015), corresponding to 23 organs and 2 cell lines. In this scenario, there is no control. Therefore, we constructed two types of TopNets independently in each tissue (Section 2.4.2). A Highest Activity TopNet (HA TopNet) where $N_i = SI$, and a z-score TopNet where $N_i = |z - \text{score}|_i$. Here N_i is the weight of node i , and SI is the normalized signal intensity (expression level). The z-score for gene i in a given tissue is computed relative to all tissues, thus using all tissues as a control for each tissue.

We assessed the tissue-specific TopNets against three gold standards (Section 2.4.3): (i) the Human Protein Atlas (HPA) (Uhlén *et al.*, 2015) where genes with ≥ 5 -fold higher abundance in each tissue are labelled tissue-specific, (ii) a set of curated tissue-relevant Disease genes (Feiglin *et al.*, 2017) and (iii) a list of genes associated with tissue-specific human diseases from the MGI (Bult *et al.*, 2019). These comparisons are carried out for 37 out of 39 tissues, as corresponding gold standards could not be obtained for the 2 cell lines. We also use a list of housekeeping genes (Eisenberg and Levanon, 2013) as a negative control. To assess the utility of the z-score TopNets, we use the same number of genes with the highest $|z - \text{score}|$ as a baseline control. Likewise, for the HA TopNets, the baseline used is the set of genes with highest expression levels.

The MGI had ≥ 25 genes with both gene expression and interaction data for 5 tissues. Of these, the HA TopNets were significantly enriched in tissue-associated genes in three tissues, and z-score TopNets in four tissues (Supplementary Table S6). In every case, the TopNet picked up equal or more gold standard genes than the corresponding baseline.

The Fisher's p -value is plotted for the overlap between the genes in the TopNets, their corresponding baselines, and gold standards Disease genes (Fig. 4, Supplementary Fig. S1) and HPA (Supplementary Figs S2, S3). Since HPA is constructed based on differential abundance, as expected, genes with top z-score are highly concordant with the HPA-derived tissue-specific genes. In all other comparisons across tissues, genes identified by PathExt agree better with gold standards than the corresponding baselines. We found 4 exceptions out of 74 comparisons (37 tissues \times 2 gold standards). Even in these cases, the pathways enriched in the TopNets are relevant to the functions of that tissue (Supplementary Information S1).

An ideal tissue-specific network should exclude housekeeping genes, which by their very definition are broadly active. We find that the TopNets identified by PathExt have this property, and are under-enriched in housekeeping genes in all but 1 case (Supplementary Table S6). This suggests that the paths prioritized by PathExt correspond to tissue-related functions rather than universally active processes.

Table 2. Comparison of PathExt with alternative approaches

Approach	PathExt		jActiveModules		DEG network						
	Active paths and the corresponding edge-induced sub-network	High scoring sub-network based on statistical significance of differential gene expression	Sub-network formed from interacting DEGs	Nodes	Edges	Gold standard	P-value				
Weightage given to Minimum number of samples to ensure optimal performance	Edges	Nodes	Nodes	Nodes	Nodes	Nodes	Nodes				
Can be used in the absence of control?	1	> 1 (enough to compute statistical significance of DEGs)	> 1 (enough to compute statistical significance of DEGs)	> 1 (enough to compute statistical significance of DEGs)	> 1 (enough to compute statistical significance of DEGs)	> 1 (enough to compute statistical significance of DEGs)	> 1 (enough to compute statistical significance of DEGs)				
Output sub-network	Yes	No	No	No	No	No	No				
Can distinguish activated and repressed processes and sub-networks?	Weighted, connected	Unweighted, connected	Unweighted, connected	Unweighted, disconnected	Unweighted, disconnected	Unweighted, disconnected	Unweighted, disconnected				
Drug	Yes	No	No	No	No	No	No				
Accession number	Edges	Edges	Edges	Edges	Edges	Edges	Edges				
Capreomycin	518	7	0.03	1066	5121	9	0.54	2426	4	0.37	
Ethionamide	394	5	0.02	1201	6941	6	0.4	161	295	6	2.8e-5
Isoniazid	401	505	2.2e-7	1145	6980	6	0.48	291	678	8	2.3e-5
Kanamycin	379	483	0.03	1107	6897	21	1.3e-5	202	358	3	0.24
Rifampicin	380	511	4.3e-6	1013	6304	7	0.44	427	926	6	0.04
Streptomycin	338	463	0.003	933	6105	16	0.001	42	45	0	0.7

Note: Response TopNets generated by PathExt are always enriched with gold standard genes, whereas existing DEG-based methods have inconsistent performance.

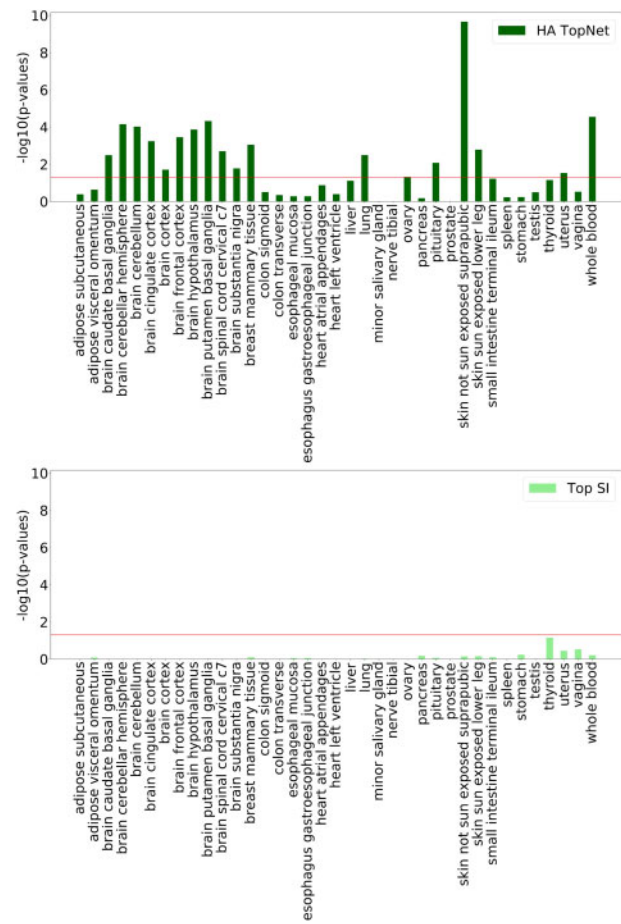


Fig. 4. Overlap with Disease genes. The p -values for overlap between the gold standard Disease genes, and nodes from HA TopNets (top) and its corresponding baseline control, Top SI (genes with highest expression) (bottom) are plotted here. Horizontal line corresponds to p -value = 0.05. The corresponding figure for z -score TopNets and top z -score can be found in [Supplementary Figure S1](#), with further discussion in [Supplementary Information S1](#). The overlap between TopNet nodes and the gold standard is better than the corresponding baseline

3.2.1 PathExt-identified pathways enriched exclusively in a tissue correspond to known tissue-relevant functions

Figure 5a shows the top pathways exclusively enriched in the HA TopNet of selected tissues (Section 2.4.4), along with literature evidence supporting each pathway-tissue association. Corresponding results for all tissues are provided in [Supplementary Tables S7](#) (HA TopNet) and [S8](#) (z -score TopNet). Some of the pathway-tissue pairs correspond to well-established functions of the tissue, such as *regulation of bile acid metabolic process* in liver ([Chiang, 2013](#)), and *ethanol catabolic process* in lung ([Bernstein, 1982](#)). PathExt reveals a few seemingly counter-intuitive associations as well. For example, *sensory perception of smell* is the top pathway exclusively enriched in the testis. However, prevalence of olfactory receptors in the testis and sperm has been experimentally verified, and testicular olfactory receptor signaling has been implicated in sperm flagellar motility ([Kang and Koo, 2012](#)). As another example, *regulation of rhodopsin mediated signaling pathway* is enriched exclusively in the pancreas. Interestingly, rhodopsin regulates insulin receptor signaling in rod photoreceptor neurons ([Rajala and Anderson, 2010](#)), and loss of Arf4, a GTPase important for localizing rhodopsin to the eye and kidney, has been shown to result in damage of exocrine pancreas in mice ([Pearring et al., 2017](#)). This surprising link between rhodopsin and the pancreas is not picked up by any of the gold standards or the controls.

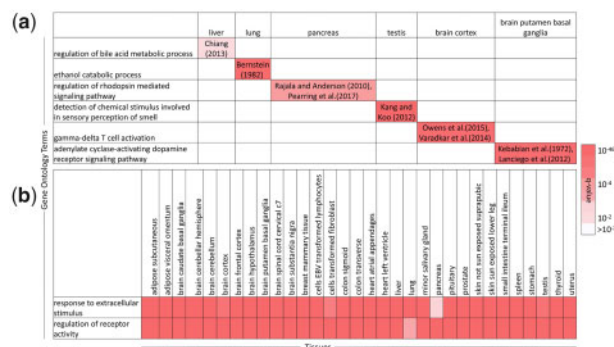


Fig. 5. Tissue-exclusive and ubiquitous pathways, HA TopNets. (a) Top GO Biological Process exclusively enriched in selected tissues, with literature evidence supporting each pathway-tissue association; (b) two processes expected to be ubiquitous, enriched in all tissues. These results suggest that PathExt can identify tissue-relevant processes despite the absence of a clear control. Corresponding result for all tissues in [Supplementary Table S7](#)

Figure 5a also highlights the specificity of functions of the different regions of the brain. For instance, *gamma-delta T cell activation* is enriched in the brain cortex. Gamma-delta T cells have been implicated in Rasmussen encephalitis, a disease characterizing inflammation of the cerebral cortex (Owens *et al.*, 2015; Varadkar *et al.*, 2014). The *adenylate cyclase-activating dopamine receptor signaling pathway* is enriched exclusively in the putamen basal ganglia region of the brain. The dorsal region of the basal ganglia comprises of the putamen, and the caudate nucleus (Lanciego *et al.*, 2012). Experiments involving homogenates of the caudate nucleus of the rat brain point at dopamine-sensitive adenylate cyclase as the receptor for dopamine in the mammalian brain (Kebabian *et al.*, 1972). This finding could indicate the presence of caudate nucleus cells in the putamen sample, or a shared function between these two adjacent regions of the brain. Several processes expected to be ubiquitous, such as *regulation of receptor activity* and *response to extracellular stimulus*, are enriched in all the tissues under consideration (Fig. 5b).

Overall, PathExt-identified tissue-specific TopNets recapitulate gold standard genes with known tissue-specific functions, and provide unique insights into tissue functions, not reflected in conventional differential expression-based analyses.

4 Discussion

We provide PathExt, a computational tool to identify sub-networks of an omics-integrated biological network, which capture the response to a perturbation, or the active processes in a particular condition. PathExt builds on our prior work which mined omics-integrated networks to (i) identify tuberculosis biomarkers (Sambarey *et al.*, 2017b), (ii) discriminate between primary and metastatic melanoma (Metri *et al.*, 2017) and (iii) identify influential genes in the condition under study (Sambaturu *et al.*, 2016). Substantially extending our prior work, PathExt provides a general framework to address all the above questions, while employing rigorous statistical significance estimation to identify critical paths. Importantly, PathExt is designed to operate even when a single sample is available for each condition, and in the absence of an experimental control sample.

In contrast to most DEG-driven methods, PathExt assigns weights to the interactions in the biological network as a function of the given omics data, thus transferring importance from individual genes to paths, and potentially capturing the way in which biological phenotypes emerge from interconnected processes. Interestingly, even though connectedness is not used as a criterion to identify sub-networks, the TopNet resulting from the identified paths forms a well-connected graph.

While the paths identified by PathExt may not constitute a comprehensive or exhaustive listing of all the active, altered processes in the system, the resulting TopNet can be thought of as a starting

point from which hypotheses can be generated. In this work, we have gathered, for each drug and each tissue, the top central genes, along with their fold change for drug exposure ([Supplementary Tables S4 and S5](#)), and z-score for human tissues ([Supplementary Tables S10 and S11](#)). Further examining the network or genomic neighborhood of these and other genes comprising the TopNet can provide additional insights, or strengthen the insights gained.

PathExt relies on two user defined parameters, the threshold k used to select the top k shortest paths, and the q-value for statistical significance of the paths selected to construct the TopNet. These values have been set at very stringent values in this paper, allowing us to focus on the most active paths. Different thresholds can give different layers of information, with different levels of false discovery.

PathExt can facilitate the study of activated and repressed processes either separately or as a whole, greatly aiding in the interpretation of the results. This is demonstrated by the INH case study, where the Activated Response TopNet highlighted the up-regulation of fatty acid synthesis, and the Repressed Response TopNet pointed to the down-regulation of its degradation. These two processes in concert suggest a possible resistance response against INH exposure, whose mechanism of action is the inhibition of fatty acid synthesis. PathExt also provides three weighting schemes, allowing users to switch between focusing exclusively on alterations, or studying both highly active and highly altered processes. Together with the ability to work with a single sample and in the absence of control, these features set PathExt apart from prior methods, providing users with a powerful, flexible, general-purpose tool for mining omics-integrated biological networks.

Funding

This work was supported by the Department of Biotechnology (DBT)—Indian Institute of Science (IISc) Partnership Program—Phase II [BT/PR27952/IN/22/212/2018]. S.H. was supported in part by National Science Foundation (NSF) award 1564785 and in part by the Intramural Research Program of the National Cancer Institute, Center for Cancer Research, National Institutes of Health (NIH).

Conflict of Interest: none declared.

References

- 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 (Methodological)*, **57**, 289–300.
- Bernstein, J. (1982) The role of the lung in the metabolism of ethanol. *Res. Commun. Chem. Pathol. Pharmacol.*, **38**, 43–56.
- Bindea, G. *et al.* (2009) ClueGO: a cytoscape plug-in to decipher functionally grouped gene ontology and pathway annotation networks. *Bioinformatics*, **25**, 1091–1093.
- Boshoff, H.I. *et al.* (2003) DnaE2 polymerase contributes to in vivo survival and the emergence of drug resistance in *Mycobacterium tuberculosis*. *Cell*, **113**, 183–193.
- Bult, C.J. *et al.*; The Mouse Genome Database Group. (2019) Mouse genome database (MGD) 2019. *Nucleic Acids Res.*, **47**, D801–D806.
- Cabusora, L. *et al.* (2005) Differential network expression during drug and stress response. *Bioinformatics*, **21**, 2898–2905.
- Carithers, L. J., and Moore, H. M. (2015) The Genotype-Tissue Expression (GTEx) Project. *Biopreservation and Biobanking*, **13**, 307–308. [10.1089/bio.2015.29031.hmm](https://doi.org/10.1089/bio.2015.29031.hmm)
- Chiang, J.Y. (2013) Bile acid metabolism and signaling. *Compr. Physiol.*, **3**, 1191–1212.
- DeJesus, M.A. *et al.* (2017) Comprehensive essentiality analysis of the *Mycobacterium tuberculosis* genome via saturating transposon mutagenesis. *MBio*, **8**, e02133.
- Dijkstra, E.W. (1959) A note on two problems in connexion with graphs. *Numer. Math.*, **1**, 269–271.
- Eisenberg, E. and Levanon, E.Y. (2013) Human housekeeping genes, revisited. *Trends Genet.*, **29**, 569–574.
- Eltringham, I. *et al.* (1999) Evaluation of reverse transcription-PCR and a bacteriophage-based assay for rapid phenotypic detection of rifampin

- resistance in clinical isolates of *Mycobacterium tuberculosis*. *J. Clin. Microbiol.*, **37**, 3524–3527.
- Esteve-Codina, A. (2018) RNA-seq data analysis, applications and challenges. *Data Anal. Omic Sci. Methods Appl.*, **82**, 71.
- Feiglin, A. et al. (2017) Comprehensive analysis of tissue-wide gene expression and phenotype data reveals tissues affected in rare genetic disorders. *Cell Syst.*, **5**, 140–148.
- Hamosh, A. et al. (2004) Online Mendelian inheritance in man (omim), a knowledgebase of human genes and genetic disorders. *Nucleic Acids Res.*, **33**, D514–D517.
- He, B. et al. (2011) A comprehensive analysis of the dynamic biological networks in HCV induced hepatocarcinogenesis. *PLoS One*, **6**, e18516.
- Howard, N.C. et al. (2018) *Mycobacterium tuberculosis* carrying a rifampicin drug resistance mutation reprograms macrophage metabolism through cell wall lipid changes. *Nat. Microbiol.*, **3**, 1099–1108.
- Ideker, T. et al. (2002) Discovering regulatory and signalling circuits in molecular interaction networks. *Bioinformatics*, **18**, S233–S240.
- Jain, P. et al. (2016) Dual-reporter mycobacteriophages (ϕ 2drms) reveal pre-existing *Mycobacterium tuberculosis* persistent cells in human sputum. *MBio*, **7**, e01023.
- Jiang, Z. et al. (2015) Whole transcriptome analysis with sequencing: methods, challenges and potential solutions. *Cell. Mol. Life Sci.*, **72**, 3425–3439.
- Kang, N. and Koo, J. (2012) Olfactory receptors in non-chemosensory tissues. *BMB Rep.*, **45**, 612–622.
- Kapopoulou, A. et al. (2011) The MycoBrowser portal: a comprehensive and manually annotated resource for mycobacterial genomes. *Tuberculosis (Edinburgh, Scotland)*, **91**, 8–13.
- Kebabian, J.W. et al. (1972) Dopamine-sensitive adenylate cyclase in caudate nucleus of rat brain, and its similarity to the “dopamine receptor”. *Proc. Natl. Acad. Sci. USA*, **69**, 2145–2149.
- Köhler, S. et al. (2014) The human phenotype ontology project: linking molecular biology and disease through phenotype data. *Nucleic Acids Res.*, **42**, D966–D974.
- Lanciego, J.L. et al. (2012) Functional neuroanatomy of the basal ganglia. *Cold Spring Harbor Perspect. Med.*, **2**, a009621–a009621.
- Ma, S. et al. (2015) Integrated modeling of gene regulatory and metabolic networks in *Mycobacterium tuberculosis*. *PLoS Comput. Biol.*, **11**, e1004543.
- McGrath, M. et al. (2014) Mutation rate and the emergence of drug resistance in mycobacterium tuberculosis. *J. Antimicrob. Chemother.*, **69**, 292–302.
- Metri, R. et al. (2017) Identification of a gene signature for discriminating metastatic from primary melanoma using a molecular interaction network approach. *Sci. Rep.*, **7**, 17314.
- Miesel, L. et al. (1998) Mechanisms for isoniazid action and resistance. In: *Genetics and Tuberculosis: Novartis Foundation Symposium 217, volume 217, pages 209–221*, Wiley Online Library, pp. 209–221.
- Milo, R. et al. (2002) Network motifs: simple building blocks of complex networks. *Science*, **298**, 824–827.
- Mishra, S. et al. (2017) Efficacy of β -lactam/ β -lactamase inhibitor combination is linked to whib4-mediated changes in redox physiology of mycobacterium tuberculosis. *Elife*, **6**, e25624.
- Mitra, K. et al. (2013) Integrative approaches for finding modular structure in biological networks. *Nat. Rev. Genet.*, **14**, 719–732.
- Owens, G.C. et al. (2015) Evidence for the involvement of gamma delta t cells in the immune response in rasmussen encephalitis. *J. Neuroinflammation*, **12**, 134.
- Pearring, J.N. et al. (2017) Loss of arf4 causes severe degeneration of the exocrine pancreas but not cystic kidney disease or retinal degeneration. *PLoS Genet.*, **13**, e1006740.
- Rajala, R.V. and Anderson, R.E. (2010) Rhodopsin-regulated insulin receptor signaling pathway in rod photoreceptor neurons. *Mol. Neurobiol.*, **42**, 39–47.
- Ritchie, M.E. et al. (2015) limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.*, **43**, e47–e47.
- Sambarey, A. et al. (2017a) Meta-analysis of host response networks identifies a common core in tuberculosis. *NPJ Syst. Biol. Appl.*, **3**, 4.
- Sambarey, A. et al. (2017b) Unbiased identification of blood-based biomarkers for pulmonary tuberculosis by modeling and mining molecular interaction networks. *EBioMedicine*, **15**, 112–126.
- Sambaturu, N. et al. (2016) Epitracer—an algorithm for identifying epicenters in condition-specific biological networks. *BMC Genomics*, **17**, 543.
- Shannon, P. et al. (2003) Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.*, **13**, 2498–2504.
- Stretch, C. et al. (2013) Effects of sample size on differential gene expression, rank order and prediction accuracy of a gene signature. *PLoS One*, **8**, e65380.
- Takayama, K. et al. (2005) Pathway to synthesis and processing of mycolic acids in mycobacterium tuberculosis. *Clin. Microbiol. Rev.*, **18**, 81–101.
- Uhlén, M. et al. (2015) Tissue-based map of the human proteome. *Science*, **347**, 1260419–1260419.
- Varadkar, S. et al. (2014) Rasmussen’s encephalitis: clinical features, pathobiology, and treatment advances. *Lancet Neurol.*, **13**, 195–205.
- Wishart, D.S. et al. (2018) Drugbank 5.0: a major update to the drugbank database for 2018. *Nucleic Acids Res.*, **46**, D1074–D1082.
- Yu, G. et al. (2012) clusterprofiler: an R package for comparing biological themes among gene clusters. *Omic J. Integr. Biol.*, **16**, 284–287.
- Zhu, J.-H. et al. (2018) Rifampicin can induce antibiotic tolerance in mycobacteria via paradoxical changes in rpoB transcription. *Nat. Commun.*, **9**, 4218.