

RESEARCH

Open Access



# Personalized single-cell networks: a framework to predict the response of any gene to any drug for any patient

Haripriya Harikumar<sup>1,2\*</sup> , Thomas P. Quinn<sup>1\*†</sup>, Santu Rana<sup>1</sup>, Sunil Gupta<sup>1</sup> and Svetha Venkatesh<sup>1</sup>

\*Correspondence:

[h.harikumar@deakin.edu.au](mailto:h.harikumar@deakin.edu.au);  
[contacttomquinn@gmail.com](mailto:contacttomquinn@gmail.com)

†Haripriya Harikumar and Thomas P. Quinn contributed equally to this work.

<sup>1</sup>Applied Artificial Intelligence Institute, Deakin University, Geelong, Australia

<sup>2</sup>Institute for Health Transformation, Deakin University, Geelong, Australia

## Abstract

**Background:** The last decade has seen a major increase in the availability of genomic data. This includes expert-curated databases that describe the biological activity of genes, as well as high-throughput assays that measure gene expression in bulk tissue and single cells. Integrating these heterogeneous data sources can generate new hypotheses about biological systems. Our primary objective is to combine population-level drug-response data with patient-level single-cell expression data to predict how any gene will respond to any drug for any patient.

**Methods:** We take 2 approaches to benchmarking a “dual-channel” random walk with restart (RWR) for data integration. First, we evaluate how well RWR can predict known gene functions from single-cell gene co-expression networks. Second, we evaluate how well RWR can predict known drug responses from individual cell networks. We then present two exploratory applications. In the first application, we combine the Gene Ontology database with glioblastoma single cells from 5 individual patients to identify genes whose functions differ between cancers. In the second application, we combine the LINCS drug-response database with the same glioblastoma data to identify genes that may exhibit patient-specific drug responses.

**Conclusions:** Our manuscript introduces two innovations to the integration of heterogeneous biological data. First, we use a “dual-channel” method to predict up-regulation and down-regulation separately. Second, we use individualized single-cell gene co-expression networks to make personalized predictions. These innovations let us predict gene function and drug response for individual patients. Taken together, our work shows promise that single-cell co-expression data could be combined in heterogeneous information networks to facilitate precision medicine.



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

## Introduction

Advances in high-throughput RNA-sequencing (RNA-Seq) have made it possible to quantify RNA presence in any biological sample [1], producing a gene expression signature that can serve as a biomarker for disease prediction [2–4] or surveillance [5, 6]. Over the last few years, single-cell RNA-Seq has risen in popularity [7]. Compared with conventional bulk RNA-Seq, which measures the average gene expression for an individual sample, single-cell RNA-Seq (scRNA-Seq) measures gene expression for an individual cell. This new mode of data makes it possible to explore tissue heterogeneity, notably tumor heterogeneity [8], by producing multiple data points per individual (i.e., one for each cell). Since genes are often understood to work in cooperative modules, the analysis of *gene co-expression networks* is commonplace. For bulk RNA-Seq, a gene co-expression network describes how genes co-occur for a population of individual samples. For scRNA-Seq, the network describes gene co-expression for a population of single cells. When these cells belong to an individual patient, the scRNA-Seq network is a kind of **personalized network** that one could use for precision medicine tasks.

Gene co-expression networks can be integrated with outside information to combine **general knowledge** (in the form of a relational database like Gene Ontology [9]) with **specific knowledge** about a sample (in the form of a co-expression network). For example, weighted gene co-expression network analysis is a popular method for functionally characterizing parts of the network, or the network as a whole [10, 11]. Although these coarse descriptions are useful, one could also combine general- and specific knowledge to make finer-level predictions about the behavior of *individual genes*. By representing each modality as a graph, multiple data streams can be combined into a **heterogeneous information network**, and then analyzed under a unified framework based on the principle of “guilt-by-association” [12] (e.g., if “a” is connected to “b” and “b” is connected to “c”, then “a” is probably connected to “c”). When the general knowledge is **gene-annotation** associations, we can (a) impute the function for genes with no known role or (b) select the most important known function. When the general knowledge is **gene-drug** response, we can predict the response of any gene to any drug. Since these inferences are tailored to the co-expression network used, they can be made personalized by using the single-cell network of an individual patient.

Random walk (RW) is a popular method that offers a general solution to the analysis of heterogeneous information networks [12, 13]. There are many variants to RW, including random walk with restart (RWR), where each step has a probability of restarting from the starting node (or a neighbor of the starting node) [14]. RW and RWR are often used in recommendation systems [15–17], but can also perform other machine learning tasks like image segmentation [18, 19], image captioning [20], or community detection [21, 22]. One advantage of RW is that it can handle missing data [23], making it a good choice for processing sparse gene annotation databases. RW and RWR have both found use in biology to find associations between genes and another data modality. For example, the “InfAcrOnt” method used an RW-based method to infer similarities between ontology terms by integrating annotations with a gene-gene interaction network [24]. Similarly, the “RWLPAP” method used RW to find lncRNA-protein associations [25], while others have used RW to predict gene-disease associations [26]. Meanwhile, RWR has been used to identify epigenetic factors within the genome [27], key genes involved in colorectal cancer [28], novel microRNA-disease associations [29], infection-related genes [30], disease-related

genes [31], and functional similarities between genes [32]. Bi-random walk, another random walk variant, has been used to rank disease genes from a protein-protein interaction network [33].

In contrast to the previous works, which make use of population-level graphs, we apply RWR to patient-level graphs, allowing us to make predictions about gene behavior that are personalized to each patient. We take 2 approaches to benchmarking “dual-channel” RWR for data integration. First, we evaluate how well RWR can predict known gene functions from single-cell gene co-expression networks. Second, we evaluate how well RWR can predict known drug responses from individual cell networks. We then present two exploratory applications. In the first application, we combine the Gene Ontology database with glioblastoma single cells from 5 individual patients to identify genes whose functions differ between cancers. In the second application, we combine the LINCS drug-response database with the same glioblastoma data to identify genes that may exhibit patient-specific drug responses. Taken together, our work shows promise that single-cell co-expression data could be combined in heterogeneous information networks to facilitate precision medicine.

## Methods

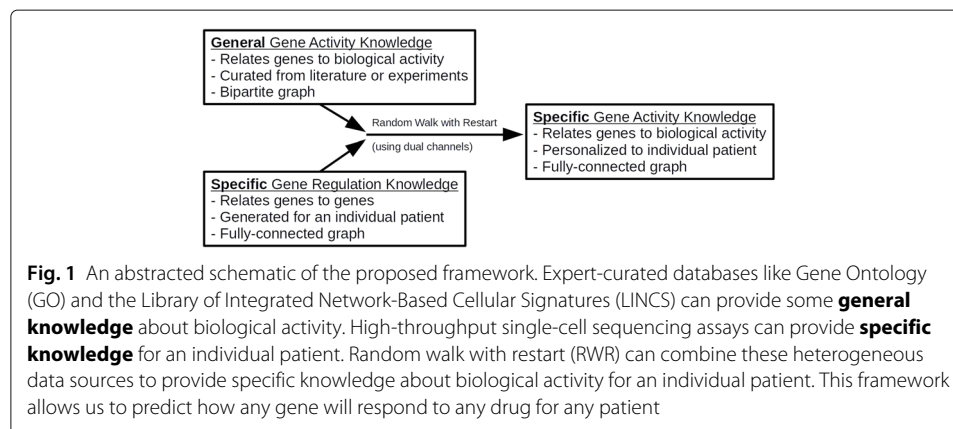
### Overview

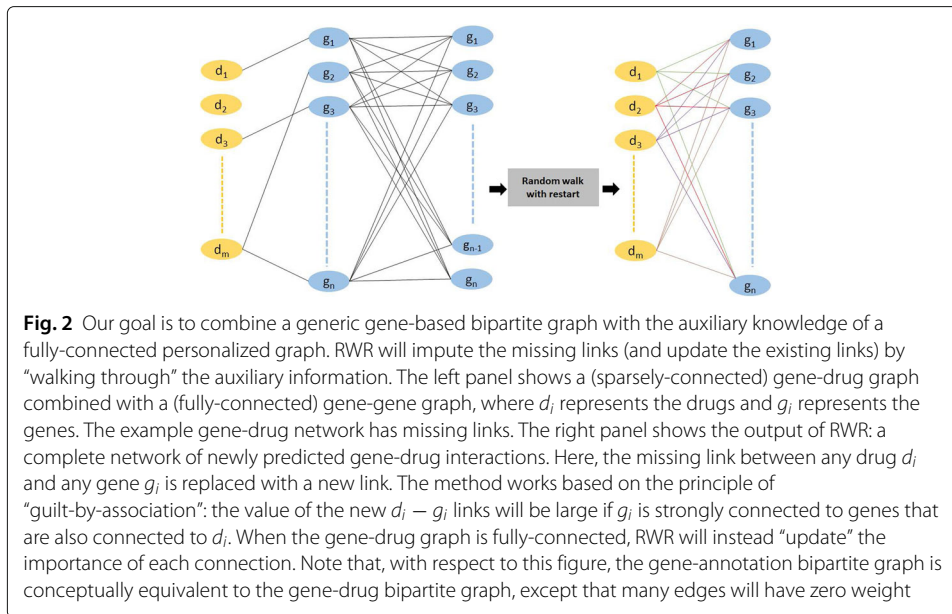
In the medical domain, gene expression can be used as a biomarker to measure the functional state of a cell. One way in which drugs mediate their therapeutic or toxic effects is by altering gene expression. However, the assays needed to test how gene expression changes in response to a drug can be expensive and time consuming. Imputation has the potential to accelerate research by “recommending” novel gene-drug responses. Random walk methods can combine sparse heterogeneous graphs based on the principle of “guilt-by-association” [12]. Figure 1 provides an abstracted schematic of the proposed framework. Figure 2 provides a visualization of the input and output for the random walk with restart (RWR) method. Figure 3 presents a bird’s-eye view of the data processing, validation, and application steps performed in this study.

### Data acquisition

The gene expression data come from two primary sources.

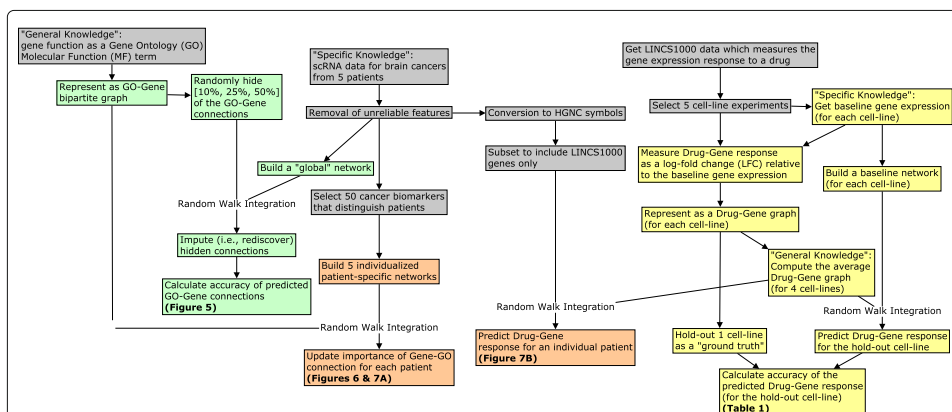
First, we acquired single-cell RNA-Seq (scRNA-Seq) expression data for 5 glioblastoma multiforme tumors [34] using the recount2 package for the R programming language [35]





(ID: SRP042161). Since scRNA-Seq data are incredibly sparse, and since the random walk with restart algorithm is computationally expensive, we elected to remove genes that had zero values in more than 25% of cells. After pre-processing, our data contain 3022 gene features and 676 single cells. These cells belong to 5 patients, with 192, 97, 97, 193, and 97 cells per patient, respectively. Finally, we randomly split the cells into 5-folds per patient so that we could estimate the variability of our downstream analyses.

Second, we acquired gene expression data from the Library of Integrated Network-Based Cellular Signatures (LINCS) [36] using the Gene Expression Omnibus (GEO) [37] (ID: GSE70138). We split these LINCS data into smaller data sets based on the cell line ID under study. We a priori included the A375 (skin; malignant melanoma), HA1E (kidney;



embryonic), HT29 (colon; adenocarcinoma), MCF7 (breast; adenocarcinoma), and PC3 (prostate; adenocarcinoma) cell lines because they were treated with the largest number of drugs.

### Defining the gene co-expression network graphs

Although correlation is a popular choice for measuring gene co-expression, correlations can yield spurious results for next-generation sequencing data [38]. Instead, we calculate the proportionality between genes using the  $\phi_s$  metric from the `propr` package for the R programming language [39]. Although this does not offer a perfect solution [40], studying gene-gene proportionality has a strong theoretical justification [38] and empirically outperforms other metrics of association for scRNA-Seq [41].

The proportionality metric describes the dissimilarity between any two genes, and ranges from  $[0, \infty)$ , where 0 indicates a perfect association. We converted this to a similarity measure  $\phi_i$  that ranges from  $[0, 1]$  by max-scaling  $\phi_i = (\max(\phi_s) - \phi_s) / \max(\phi_s)$ , such that  $\phi_i = 1$  when  $\phi_s = 0$ . A gene-gene matrix of  $\phi_i$  scores is analogous to a gene-gene matrix of correlation coefficients, and constitutes our gene co-expression network. We calculated the  $\phi_i$  co-expression network for the entire scRNA-Seq data set (1 network), for 5 folds of 5 patients (25 networks total), and for each of the 5 LINCS drug-free cell lines (5 networks total). All co-expression networks are available from <https://zenodo.org/record/3522494>.

### Defining the bipartite graphs

Consider a graph  $G$  with  $V = 1..N$  vertices, with positive and negative edges. The graphs used for our analyses are composed for two parts: a (general knowledge) bipartite graph and a (specific knowledge) fully-connected gene co-expression graph. For a bipartite graph, the vertex set  $V$  can be separated into two distinct sets,  $V_1$  and  $V_2$ , such that no edges exist within either set. For a fully-connected (or complete) graph, there exists an edge between every pair of vertices within one set. For the graph  $G$ , the bipartite and fully-connected graphs are joined via the common vertex set  $V_1$  that contains genes and  $V_2$  contains annotations or drugs.

We constructed two types of bipartite graphs: the **gene-annotation graph** and the **gene-drug graph**. First, we made the gene-annotation graph from the Gene Ontology Biological Process database [9] via the `AnnotationDbi` and `org.Hs.eg.db` Bioconductor packages. An edge exists whenever a gene is associated with an annotation. Second, we made the gene-drug graphs using the LINCS data. For each cell line, we computed a gene-drug graph by calculating the log-fold change between the median of the drug-treated cell's expression and the median of the drug-naive cell's expression. This results in a fully-connected and weighted bipartite graph, where a large positive value means that the drug causes the gene to up-regulate (and *vice versa*). All bipartite graphs are available from <https://zenodo.org/record/3522494>.

### Dual-channel random walk with restart (RWR)

Traditional RWR methods can only perform a random walk on graphs with positive edge weights [13]. Since the response of a gene to a drug is directional (up-regulated or down-regulated), we chose to use a modified RWR method, proposed by [42], that handles graphs with both positive and negative edge weights. Random walk requires tran-

sition probability matrices to decide the next step in the walk. The Chen et al. transition probability matrices can be computed based on the following equations:

$$P(x_j|x_i) = \frac{|e_{ij}|}{\sum_{l \in N(x_i)} |e_{il}|} \tag{1}$$

$$P(x_j^+|x_i^+) = P(x_j^-|x_i^-) = \begin{cases} P(x_j|x_i), & \text{if } e_{ij} \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

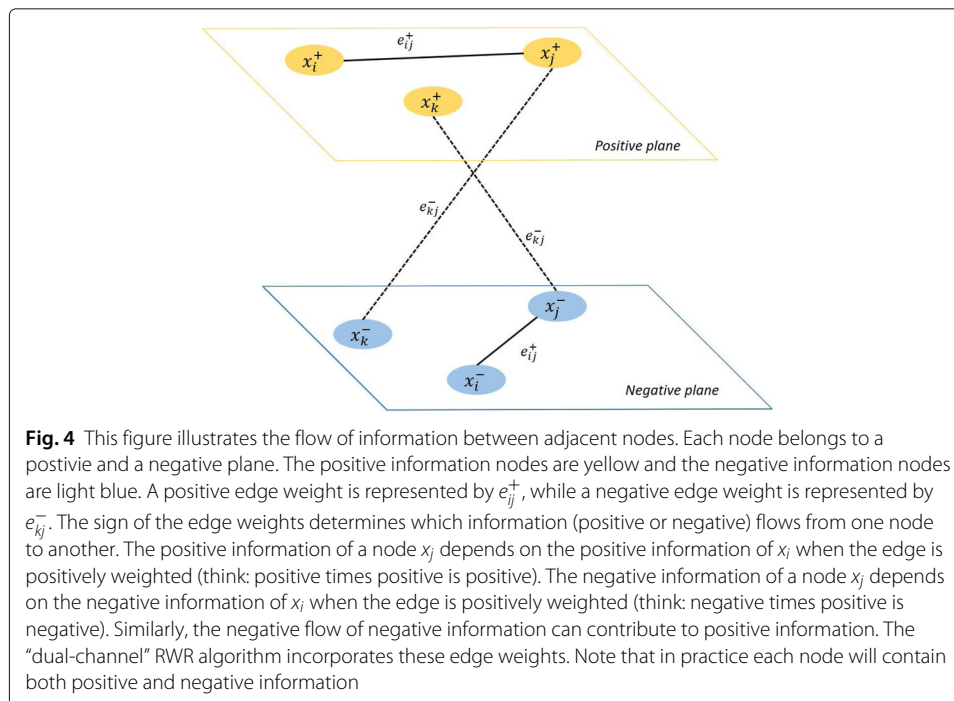
$$P(x_j^-|x_i^+) = P(x_j^+|x_i^-) = \begin{cases} P(x_j|x_i), & \text{if } e_{ij} < 0 \\ 0, & \text{otherwise} \end{cases}$$

These equations separate out the positive (and negative) transitions, and are used to calculate the total positive (and negative) information flow for each node. They are fixed for all steps.

Though the transition probabilities are computed separately, the information accumulated in a node depends on both the positive and negative information which flows through the node. For example, the positive information in a node depends on the negative information of any neighboring node connected by a negative edge weight (think: negative times negative is positive). Likewise, negative information in a node depends on the positive information in a neighboring node connected by a negative edge weight, and *vice versa* (think: negative times positive is negative). Figure 4 illustrates the information flow to a node  $x_j$  from two neighbors.

The flow of information between the positive “plane” of the graph to the negative “plane” of the graph can be formulated with the equations:

$$P(x_j^+)_k = \left[ \sum_{x_i \in N(x_j) \ \& \ e_{ij} \geq 0} P(x_i^+)_{k-1} P(x_j^+ | x_i^+) \right] + \left[ \sum_{x_i \in N(x_j) \ \& \ e_{ij} < 0} P(x_i^-)_{k-1} P(x_j^+ | x_i^-) \right] \tag{2}$$



$$P(x_j^-)_k = \left[ \sum_{x_i \in N(x_j) \ \& \ e_{ij} \geq 0} P(x_i^-)_{k-1} P(x_j^- | x_i^-) \right] + \left[ \sum_{x_i \in N(x_j) \ \& \ e_{ij} < 0} P(x_i^+)_{k-1} P(x_j^- | x_i^+) \right] \tag{3}$$

where the probability  $P(x_j^+)_k$  is updated at each step  $k = 2 \dots 10000$ .

RWR always considers a probability  $\alpha$  to return back to the original nearest neighboring nodes at each step in the random walk. This is used to weigh the importance of node-specific information with respect to the whole graph, including for long walks:

$$P_{rst}(x_j^+)_k = (1 - \alpha) \times P(x_j^+)_k + \alpha \times P(x_j^+)_2 \tag{4}$$

$$P_{rst}(x_j^-)_k = (1 - \alpha) \times P(x_j^-)_k + \alpha \times P(x_j^-)_2 \tag{5}$$

where the restart probability  $P_{rst}(x_j^+)_k$  is updated at each step  $k = 2 \dots 10000$ , and  $P(x_j^+)_2$  is the probability after the first update. These equations find the positive and negative restart information with respect to the node  $x_j$ . Each  $P_{rst}(x_j^+)_k$  is a vector of probabilities that together sum to 1. This probability has two parts: the global information and the local information. The local information is the initial probability with respect to the nearest neighbors of node  $x_j$ , and is denoted by  $P(x_j^+)_2$  [or  $P(x_j^-)_2$ ] (i.e., the probability after the first update). The restart probability  $\alpha$  is chosen from the range  $[0, 1]$ , where a higher value weighs the local information more than the global information. We chose  $\alpha = 0.1$  to place a larger emphasis on the global information. This is also the value used by [42]. When applied to our synthetic data (see Additional file 1: Appendix), it produced expected results.

**Analysis of random walk with restart (RWR) scores**

For each gene, the RWR algorithm returns a vector of probabilities that together sum to 1. According to the guilt-by-association assumption, we interpret these probabilities to indicate the strength of the connection between the reference gene and each target. Since we are only interested in gene-annotation and gene-drug relationships, we exclude all gene-gene probabilities. Viewing the probability distribution as a composition (c.f., [43]), we perform a centered log-ratio transformation of each probability vector subset. This transformation normalizes the probability distributions so that we can compare them between samples [44, 45]. We define the RWR score  $r_{ga}^+$  (or  $r_{ga}^-$ ) for each gene-annotation connection as the transform of its RWR probability:

$$r_{ga}^+ = \log \frac{p_{ga}^+}{\sqrt[A]{\prod_i^A p_{gi}^+}} \tag{6}$$

$$r_{ga}^- = \log \frac{p_{ga}^-}{\sqrt[A]{\prod_i^A p_{gi}^-}} \tag{7}$$

for a bipartite graph describing  $g = 1 \dots G$  genes and  $a = 1 \dots A$  annotations (or  $A$  drugs), where  $\mathbf{p}_g^+ = P_{rst}(x_g^+) = [p_{g1}^+, \dots, p_{gA}^+]$  (i.e., from the final step). These transformed RWR scores can be used for univariate statistical analyses, such as an analysis of variance (ANOVA) (c.f., analysis of compositional data [46, 47]).

## Benchmark validation

### *Validation of gene-annotation prediction*

Our strategy to validate RWR for gene-annotation prediction involves “hiding” known functional associations and seeing whether the RWR algorithm can re-discover them. This is done by turning 1s into 0s in the bipartite graph, a process we call “sparsification”. Our sparsification procedure works in 4 steps. First, we combine the original GO BP (or MF) bipartite graph with the master single-cell co-expression graph. Second, we subset the graph to include 25% of the gene annotations and 25% of the genes (this is done to reduce the computational overhead). Third, we randomly hide [10,25,50] percent of the gene-annotation connections from the bipartite subgraph. Since this random selection could cause a feature to lose all connections, we use a constrained sampling strategy: the subsampled graph must contain at least one non-zero entry for each feature. Fourth, we apply the RWR algorithm to the sparsified and non-sparsified graphs, separately. We repeat this process 25 times, using a different random graph each time. By comparing the RWR scores between the hidden and unknown connections, we can determine whether our method rediscovers hidden connections.

### *Validation of gene-drug prediction*

We use a different strategy to validate RWR for drug-response prediction. Since we have the gene-drug and gene-gene interaction data for 5 cell lines (A375, HA1E, HT29, MCF7 and PC3), we can set aside the known gene-drug responses for 1 cell line (PC3) as a “ground truth” test set. Then, we can use a composite of the remaining 4 gene-drug graphs to predict the gene-drug responses for the withheld cell line.

This is done in two steps. First, we use the averaged gene-drug data for 4 cell lines (a general drug graph) and the gene-gene data for PC3 (a specific gene graph) to impute the gene-drug response for PC3 (a specific drug graph). In the second step, we use the gene-drug data for PC3 (a specific drug graph) and its corresponding gene-gene data (a specific gene graph) to calculate the “ground truth” RWR scores for PC (a specific drug graph). The “ground truth” is the RWR scores when all PC3 drug-response experiments have been performed. With these two outputs, we can calculate the agreement between the imputed and “ground truth” RWR scores (using Spearman’s correlation, MSE, and accuracy).

### **Exploratory application of gene-drug prediction**

Having demonstrated that RWR can perform well for single-cell co-expression networks, and can make meaningful drug-response predictions from composite LINCS data, we combine these heterogeneous data sources to make personalized drug-response predictions for individual single-cell networks. This requires some data munging. First, we transform the ENGS features used by the single-cell data into the HGNC features used by LINCS (only including genes with a 1-to-1 mapping, resulting in 181 genes). Second, we build an HGNC co-expression network with  $\phi_i$  (for 5 folds of 5 patients, yielding 25 networks total). Third, we combine the composite LINCS gene-drug bipartite graph with each of the 25 HGNC single-cell networks. Fourth, we use our RWR algorithm to predict how 181 genes would respond to 1732 drugs for each patient fold. As above, we perform an analysis of variance (ANOVA) to detect inter-patient differences.



## Results and discussion

### Why use single cells?

In this study, we analyze a previously published single-cell data set that measured the gene expression for 5 glioblastoma patients. A principal components analysis of these data show that the major axes of variance tend to group the cells according to the patient-of-origin. Indeed, an ANOVA of gene expression with respect to patient ID reveals that 2204 of the 3022 genes have significantly different expression in at least one patient (FDR-adjusted  $p < .05$ ). This suggests that single-cell gene expression is unique to each patient.

Although it is possible to obtain sample-specific *gene expression* using bulk RNA-Seq, our approach to data integration exploits the graphical structure of sample-specific *gene co-expression networks*. scRNA-Seq, generating multiple measurements per individual, makes the computation of sample-specific co-expression networks straightforward (though others have proposed ways to estimate these from bulk RNA-Seq [48, 49]).

### Validation of gene-annotation prediction

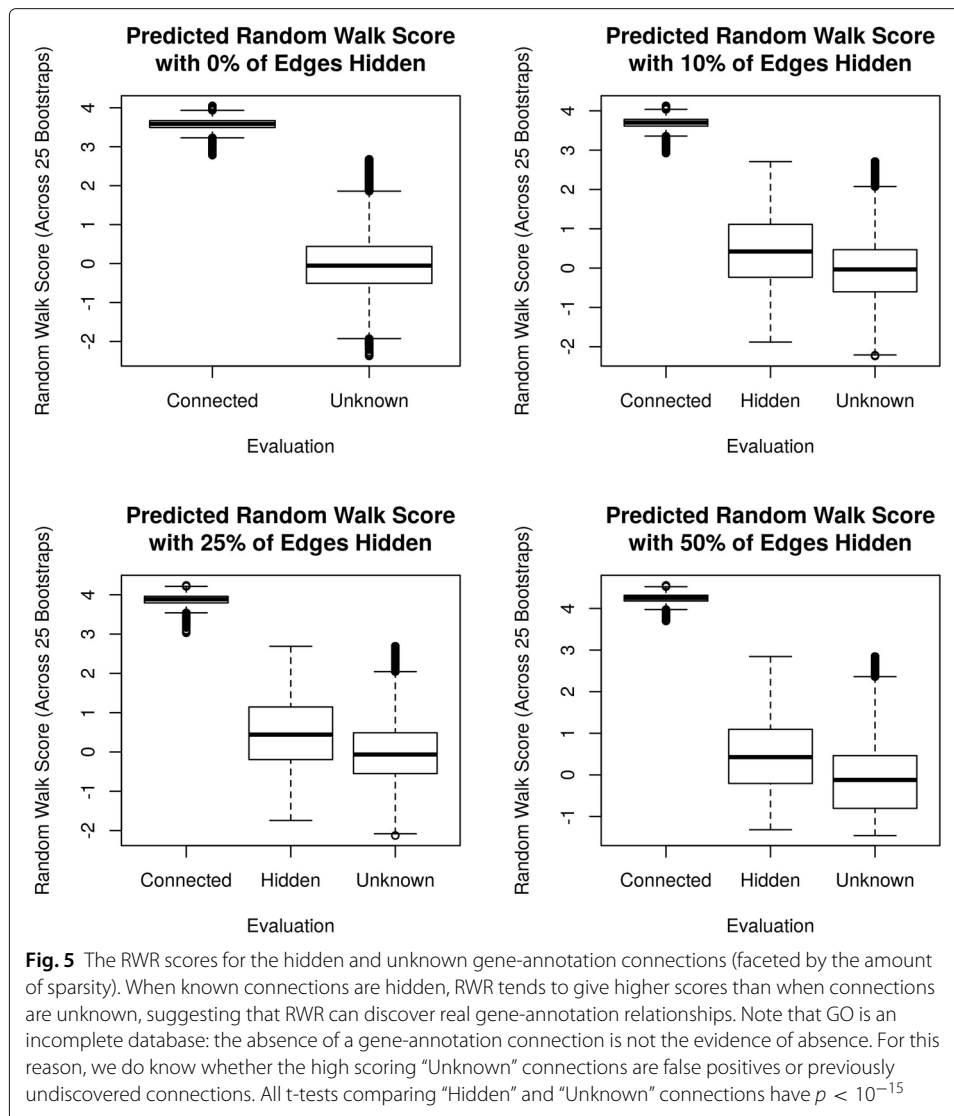
The Gene Ontology (GO) project has curated a database which relates genes to biological processes (BP) and molecular functions (MF) (called **annotations**). The GO database has widespread use in bioinformatics for assigning “functional” relevance to sets of gene biomarkers [50]. Although GO organizes the semantic relationships between annotations as a directed acyclic graph, we could more simply represent the relationships as a bipartite graph. By combining a (fully-connected) gene co-expression graph with a (sparsely-connected) gene-annotation bipartite graph, RWR can predict new gene-annotation connections.

To test whether the RWR predictions are meaningful, we “hid” a percentage of known gene-annotation links (by turning 1s into 0s in the bipartite graph), and compared the RWR scores for the *hidden* gene-annotation links against those for the *unknown* links (see [Methods](#) for a definition of the RWR score). Figure 5 shows that the RWR scores for hidden connections are appreciably larger than for the unknown connections, confirming that RWR can discover real gene-annotation relationships from a single-cell gene co-expression network.

### Exploratory application of gene-annotation prediction

Since single-cell RNA-Seq assays measure RNA for multiple cells per patient, we can use these data to build a personalized graph that describes the gene-gene relationships for an individual patient. In order to estimate the variation in these personalized graphs, we divided the cells from each sample into 5 folds (giving us 5 networks per-patient). Above, we show that RWR can discover real gene-annotation relationships. By combining the personalized graph (a kind of *specific knowledge*) with a gene-annotation bipartite graph (a kind of *general knowledge*), the RWR algorithm will score the gene-annotation connections for a given patient. From this, we can identify genes that may have a different functional importance in one cancer versus the others.

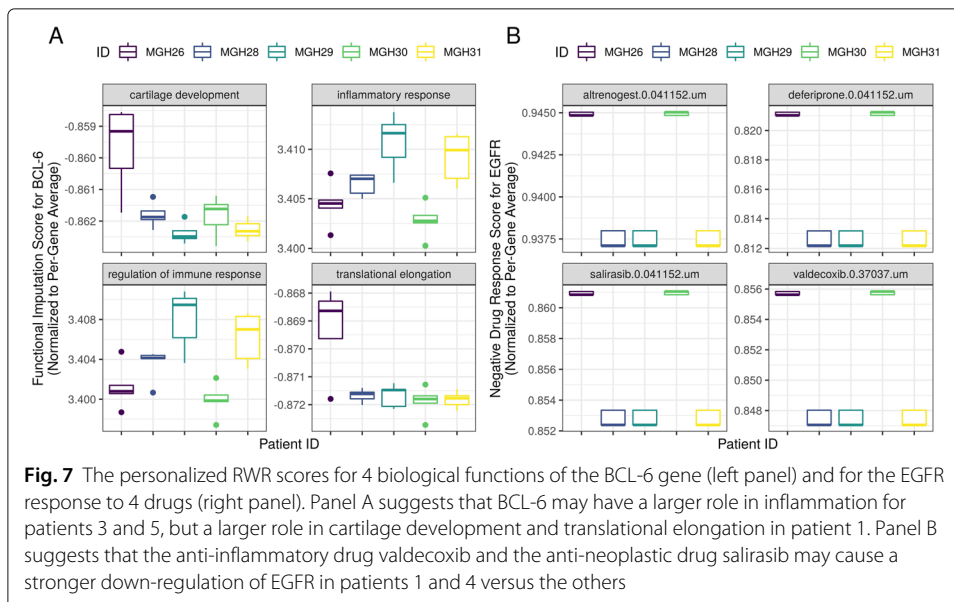
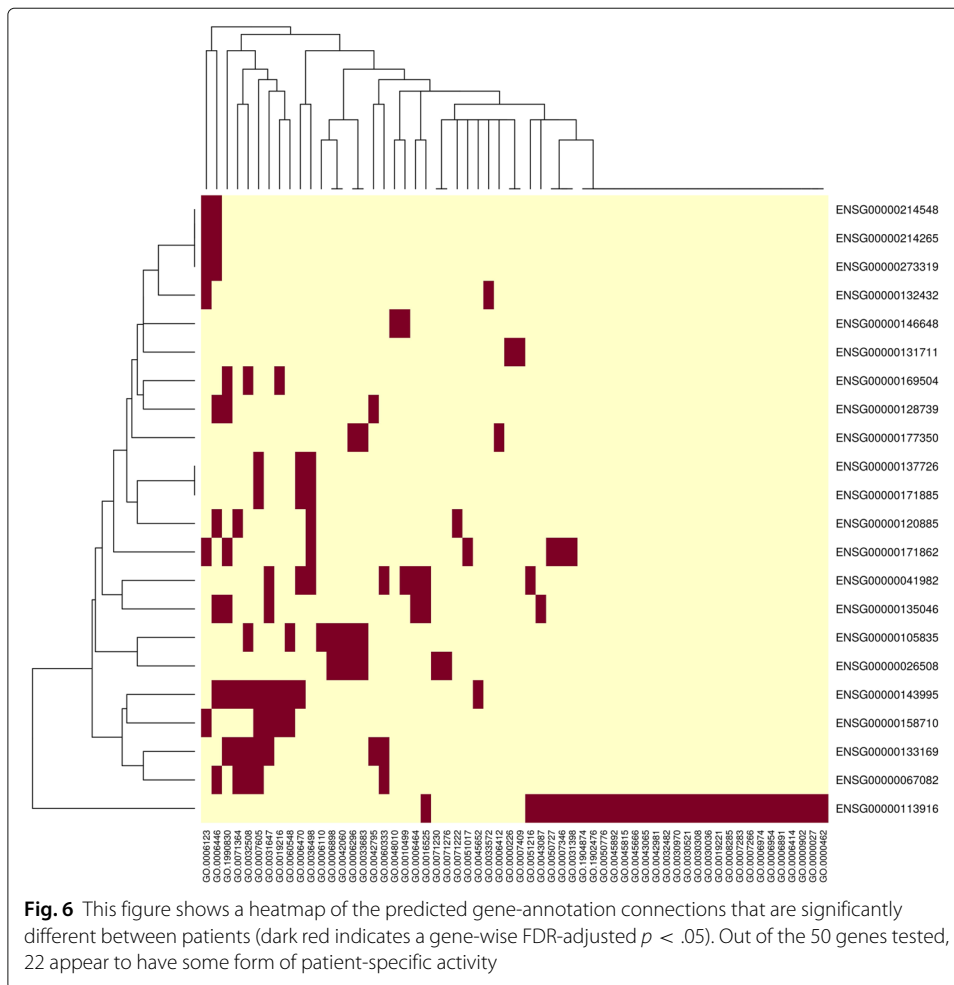
Taking a subset of the 50 genes with the largest inter-patient differences, we use RWR to compute personalized RWR scores. This results in 25 matrices (for 5 folds of 5 patients), each with 50 rows (for genes) and 369 columns (for BP annotations). Performing an ANOVA on each gene-annotation connection results in a matrix of 50x369  $p$ -values. Figure 6 shows a heatmap of the significant gene-annotation connections (dark red indi-



cates a gene-wise FDR-adjusted  $p < .05$ ). Figure 7A plots the per-patient RWR scores for 4 annotations of the BCL-6 gene that significantly differ between patients. BCL-6 is an important biomarker whose increased expression is associated with worse outcomes in glioblastoma [51]. Our analysis suggests that BCL-6 could have a larger role in inflammation for patients 3 and 5, but a larger role in cartilage development and translational elongation in patient 1. Of course, this hypothesis requires experimental validation.

#### Validation of gene-drug prediction

The NIH LINCS program has generated a large amount of data on how the gene expression signatures of cell lines change in response to a drug. By conceptualizing the baseline (drug-free) gene co-expression network as a complete graph of *specific knowledge*, and by re-factoring the average gene-drug response as a (weighted) bipartite graph of *general knowledge*, we can apply the same RWR algorithm to predict a cell's gene expression response to any drug. Since the modified RWR algorithm contains two channels—a



positive and negative channel—we can predict up-regulation or down-regulation events separately.

To test whether RWR can make accurate predictions about how a gene in a cell would respond to a drug, we ran the RWR algorithm on the baseline (drug-free) gene co-expression graph of the PC3 cell line using a composite gene-drug graph of 4 different cell lines. We then compared these RWR scores with a “ground truth” (i.e., the RWR scores for when all PC3 drug-response experiments have been performed). The agreement between the composite gene-drug RWR scores and the “ground truth” gene-drug RWR scores tells us how well the composite gene-drug map generalizes to new cell types. Table 1 shows that agreement is high, especially for the top up-regulation and down-regulation events. This confirms that our composite gene-drug graph is useful for drug-response prediction.

### Exploratory application of gene-drug prediction

The RWR algorithm can combine specific knowledge and general knowledge from disparate sources to make personalized recommendations. This makes RWR a potentially valuable tool for precision medicine.

As an exploratory analysis, we combine the personalized gene co-expression networks with the composite gene-drug graph from LINCS. By running the RWR algorithm on these two data streams, the RWR scores will now suggest how the expression of any gene might change in response to any drug for each of the 5 glioblastoma patients. Using an ANOVA, we identify hundreds of gene-drug connections with RWR scores that differ significantly between patients (gene-wise FDR-adjusted  $p < .05$ ).

Figure 7B shows an example of drugs that have different (negative channel) RWR scores for EGFR. It suggests that the anti-inflammatory drug valdecoxib and the anti-neoplastic drug salirasib may cause a stronger down-regulation of EGFR (a pan-cancer oncogene [52]) in patients 1 and 4 versus the others. The Supplementary Information includes a complete table of the unadjusted ANOVA  $p$ -values for the gene-drug inter-patient differences available in <https://zenodo.org/record/3743897>.

### Limitations

We deployed our framework on only 5 individual patients. As such, we lack a sufficient sample size to test whether any inter-patient differences could be explained by known demographic or clinical phenotypes. It is worth noting that cancer cells are very heterogeneous and, depending on the location of the sample collection, the composition of cell types (and thus gene expression profiles) can change dramatically. As such, factors other than the patient-specific tumour profile, such as batch effects, could account for differences in the sample-specific gene co-expression networks. Such differences may be difficult to account for without careful experimental design and standardization. Although the scope of this paper is to prove the concept, we wish to remind the readers

**Table 1** Overall agreement (Spearman’s correlation and MSE) and the accuracy of the overlap (for the top 5%, 10%, 25%, and 50% predicted scores), as calculated separately for the positive and negative channels. Overall, agreement is high, especially for the top up-regulation and down-regulation events

	Correlation	MSE	Top 5% (ACC)	Top 10% (ACC)	Top 25% (ACC)	Top 50% (ACC)
Positive Channel	0.7173	0.2022	0.9279	0.8857	0.8574	0.8427
Negative Channel	0.5502	0.2176	0.9450	0.8946	0.7578	0.7053

that much care should be taken when translating the methodology to real-world clinical problem-solving.

In the absence of experimental validation, we support our analyses using 2 forms of *in silico* validation, which together demonstrate that RWR can integrate sparse heterogeneous data to discover real biological activity. Although we find the *in silico* validation encouraging, we acknowledge that RWR is merely a prediction tool that recommends hypotheses, and that these predictions may change when the source of general knowledge changes. Experimental validation is needed to determine whether these hypotheses prove true in practice. Further work is needed to validate the clinical relevance of the proposed framework.

## Conclusions

This manuscript describes a framework for combining patient-specific single-cell networks with public drug-response data to make personalized predictions about drug response. Importantly, our approach makes use of a generic framework, and so can be applied to combine many kinds of data. We think the targeted analysis of personalized single-cell networks is promising, and could offer a new direction for precision medicine research.

We conclude with some perspectives on what the future of personalized network analysis may hold. Although RWR can handle sparse heterogeneous data, the positive and negative information obtained for each node can be infinitesimally small. One might address this by first transforming the RWR probabilities. Otherwise, we note that RWR is computationally expensive, making the analysis of high-dimensional data prohibitively slow. One might address this by pre-training a deep neural network to provide an approximate RWR solution. These improvements could help scale personalized predictions to larger graphs.

## Abbreviations

RW: Random walk; RWR: Random walk with restart; RNA-Seq: RNA sequencing; MSE: Mean square error; scRNA-Seq: Single-cell RNA sequencing; LINCS: Library of Integrated Network-Based Cellular Signatures; GO: Gene Ontology; BP: Biological Process; MF: Molecular Function; ANOVA: Analysis of variance

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13040-021-00263-w>.

**Additional file 1:** Appendix.

## Acknowledgments

Not applicable.

## Authors' contributions

HH implemented the RWR algorithm and applied it to the graphical data. TPQ prepared the graph data and performed the analysis of the resultant RWR scores. HH and TPQ reviewed the literature, designed the experiments, and drafted the manuscript. All authors helped conceptualize the project and revise the manuscript. All authors read and approved the final manuscript.

## Availability of data and material

The raw data are publicly available from the resources described in the [Methods](#). All gene co-expression and bipartite graphs used in these analyses are available from <https://zenodo.org/record/3522494>.

## Declarations

### Ethics approval and consent to participate

Not applicable.

**Consent for publication**

Not applicable.

**Competing interests**

The authors declare that they have no competing interests.

Received: 21 January 2021 Accepted: 10 May 2021

Published online: 05 August 2021

**References**

1. Metzker ML. Sequencing technologies — the next generation. *Nat Rev Genet.* 2010;11(1):31–46. <https://doi.org/10.1038/nrg2626>.
2. Golub TR, Slonim DK, Tamayo P, Huard C, Gaasenbeek M, Mesirov JP, Coller H, Loh ML, Downing JR, Caligiuri MA, Bloomfield CD, Lander ES. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science (New York, N.Y.)* 1999;286(5439):531–7.
3. Alon U, Barkai N, Notterman DA, Gish K, Ybarra S, Mack D, Levine AJ. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc Natl Acad Sci U S A.* 1999;96(12):6745–50.
4. van 't Veer LJ, Dai H, van de Vijver MJ, He YD, Hart AAM, Mao M, Peterse HL, van der Kooy K, Marton MJ, Witteveen AT, Schreiber GJ, Kerkhoven RM, Roberts C, Linsley PS, Bernards R, Friend SH. Gene expression profiling predicts clinical outcome of breast cancer. *Nature.* 2002;415(6871):530–6. <https://doi.org/10.1038/415530a>.
5. Noto K, Majidi S, Edlow AG, Wick HC, Bianchi DW, Slonim DK. CSAX: Characterizing Systematic Anomalies in eXpression Data. *J Comput Biol.* 2015;22(5):402–13. <https://doi.org/10.1089/cmb.2014.0155>.
6. Quinn TP, Nguyen T, Lee SC, Venkatesh S. Cancer as a Tissue Anomaly: Classifying Tumor Transcriptomes Based Only on Healthy Data. *Front Genet.* 2019;10. <https://doi.org/10.3389/fgene.2019.00599>.
7. Gawad C, Koh W, Quake SR. Single-cell genome sequencing: current state of the science. *Nat Rev Genet.* 2016;17(3):175–88. <https://doi.org/10.1038/nrg.2015.16>.
8. Lawson DA, Kessenbrock K, Davis RT, Pervolarakis N, Werb Z. Tumour heterogeneity and metastasis at single-cell resolution. *Nat Cell Biol.* 2018;20(12):1349–60. <https://doi.org/10.1038/s41556-018-0236-7>.
9. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC, Richardson JE, Ringwald M, Rubin GM, Sherlock G. Gene Ontology: tool for the unification of biology. *Nat Genet.* 2000;25(1):25–9. <https://doi.org/10.1038/75556>.
10. Langfelder P, Horvath S. Eigengene networks for studying the relationships between co-expression modules. *BMC Syst Biol.* 2007;1:54. <https://doi.org/10.1186/1752-0509-1-54>.
11. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics.* 2008;9:559. <https://doi.org/10.1186/1471-2105-9-559>.
12. Tsuyuzaki K, Nikaido I. Biological Systems as Heterogeneous Information Networks: A Mini-review and Perspectives. 2017. <https://arxiv.org/abs/1712.08865v1>. Accessed 29 Oct 2019.
13. Pearson K. The problem of the random walk. *Nature.* 1905;72(1867):342.
14. Tong H, Faloutsos C, Pan J-Y. Fast random walk with restart and its applications. In: Sixth International Conference on Data Mining (ICDM'06). IEEE; 2006. p. 613–22.
15. Bogers T. Movie recommendation using random walks over the contextual graph. In: Proc. of the 2nd Intl. Workshop on Context-Aware Recommender Systems. Citeseer; 2010.
16. Cooper C, Lee SH, Radzik T, Siantos Y. Random walks in recommender systems: exact computation and simulations. In: Proceedings of the 23rd International Conference on World Wide Web. ACM; 2014. p. 811–6.
17. Keramarrec A-M, Leroy V, Moin A, Thraves C. Application of random walks to decentralized recommender systems. In: International Conference On Principles Of Distributed Systems. Springer; 2010. p. 48–63.
18. Grady L. Random walks for image segmentation. *IEEE Trans Patt Anal Mach Intell.* 2006;28(11):1768–83.
19. Jha SK, Bannerjee P, Banik S. Random walks based image segmentation using color space graphs. *Procedia Technol.* 2013;10:271–8.
20. Pan J-Y, Yang H-J, Faloutsos C, Duygulu P. Gcap: Graph-based automatic image captioning. In: 2004 Conference on Computer Vision and Pattern Recognition Workshop. IEEE; 2004. p. 146.
21. Pons P, Latapy M. Computing communities in large networks using random walks. In: International Symposium on Computer and Information Sciences. Springer; 2005. p. 284–93.
22. Kuncheva Z, Montana G. Community detection in multiplex networks using locally adaptive random walks. In: Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015. ACM; 2015. p. 1308–15.
23. Greenland S, Mansournia MA, Altman DG. Sparse data bias: a problem hiding in plain sight. *bmj.* 2016;352:1981.
24. Cheng L, Jiang Y, Ju H, Sun J, Peng J, Zhou M, Hu Y. Infacront: calculating cross-ontology term similarities using information flow by a random walk. *BMC genomics.* 2018;19(1):919.
25. Zhao Q, Liang D, Hu H, Ren G, Liu H. Rwlpa: Random walk for lncrna-protein associations prediction. *Protein Pept Lett.* 2018;25(9):830–7.
26. Zhao Z-Q, Han G-S, Yu Z-G, Li J. Laplacian normalization and random walk on heterogeneous networks for disease-gene prioritization. *Comput Biol Chem.* 2015;57:21–8.
27. Li J, Chen L, Wang S, Zhang Y, Kong X, Huang T, Cai Y-D. A computational method using the random walk with restart algorithm for identifying novel epigenetic factors. *Mol Gen Genomics.* 2018;293(1):293–301.
28. Cui X, Shen K, Xie Z, Liu T, Zhang H. Identification of key genes in colorectal cancer using random walk with restart. *Mol Med Rep.* 2017;15(2):867–72.
29. Sun J, Shi H, Wang Z, Zhang C, Liu L, Wang L, He W, Hao D, Liu S, Zhou M. Inferring novel lncrna–disease associations based on a random walk model of a lncrna functional similarity network. *Mol Biosyst.* 2014;10(8):2074–81.

30. Zhu L, Su F, Xu Y, Zou Q. Network-based method for mining novel hpv infection related genes using random walk with restart algorithm. *Biochim Biophys Acta (BBA)-Mol Basis Dis.* 2018;1864(6):2376–83.
31. Valdeolivas A, Tichit L, Navarro C, Perrin S, Odelin G, Levy N, Cau P, Remy E, Baudot A. Random walk with restart on multiplex and heterogeneous biological networks. *Bioinformatics.* 2018;35(3):497–505.
32. Peng J, Zhang X, Hui W, Lu J, Li Q, Liu S, Shang X. Improving the measurement of semantic similarity by combining gene ontology and co-functional network: a random walk based approach. *BMC Syst Biol.* 2018;12(2):18.
33. Xie M, Hwang T, Kuang R. Prioritizing disease genes by bi-random walk. In: *Pacific-Asia Conference on Knowledge Discovery and Data Mining.* Springer; 2012. p. 292–303.
34. Patel AP, Tirosi I, Trombetta JJ, Shalek AK, Gillespie SM, Wakimoto H, Cahill DP, Nahed BV, Curry WT, Martuza RL, Louis DN, Rozenblatt-Rosen O, Suvà ML, Regev A, Bernstein BE. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science.* 2014;344(6190):1396–401. <https://doi.org/10.1126/science.1254257>.
35. Collado-Torres L, Nellore A, Jaffe AE. recount workflow: Accessing over 70,000 human RNA-seq samples with Bioconductor. *F1000Research.* 2017;6:1558. <https://doi.org/10.12688/f1000research.12223.1>.
36. Koletti A, Terryn R, Stathias V, Chung C, Cooper DJ, Turner JP, Vidović D, Forlin M, Kelley TT, D'Urso A, Allen BK, Torre D, Jagodnik KM, Wang L, Jenkins SL, Mader C, Niu W, Fazel M, Mahi N, Pilarczyk M, Clark N, Shamsaei B, Meller J, Vasiliaskas J, Reichard J, Medvedovic M, Ma'ayan A, Pillai A, Schürer SC. Data Portal for the Library of Integrated Network-based Cellular Signatures (LINCS) program: integrated access to diverse large-scale cellular perturbation response data. *Nucleic Acids Res.* 2018;46(D1):558–66. <https://doi.org/10.1093/nar/gkx1063>.
37. Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res.* 2002;30(1):207–10. <https://doi.org/10.1093/nar/30.1.207>.
38. Lovell D, Pawlowsky-Glahn V, Egozcue JJ, Marguerat S, Bähler J. Proportionality: A Valid Alternative to Correlation for Relative Data. *PLoS Comput Biol.* 2015;11(3): <https://doi.org/10.1371/journal.pcbi.1004075>.
39. Quinn TP, Richardson MF, Lovell D, Crowley TM. propr: An R-package for Identifying Proportionally Abundant Features Using Compositional Data Analysis. *Sci Rep.* 2017;7(1):16252. <https://doi.org/10.1038/s41598-017-16520-0>.
40. Erb I, Notredame C. How should we measure proportionality on relative gene expression data?. *Theory Biosci.* 2016;135:21–36. <https://doi.org/10.1007/s12064-015-0220-8>.
41. Skinnider MA, Squair JW, Foster LJ. Evaluating measures of association for single-cell transcriptomics. *Nat Methods.* 2019;16(5):381–6. <https://doi.org/10.1038/s41592-019-0372-4>.
42. Chen Y-C, Lin Y-S, Shen Y-C, Lin S-D. A modified random walk framework for handling negative ratings and generating explanations. *ACM Trans Intell Syst Technol (tISt).* 2013;4(1):12.
43. Erb I, Ay N. The information-geometric perspective of Compositional Data Analysis. *arXiv preprint arXiv:2005.11510.* 2020.
44. Boogaart K. G. v. d., Tolosana-Delgado R. *Fundamental Concepts of Compositional Data Analysis.* In: *Analyzing Compositional Data With R. Use R!* Springer; 2013. p. 13–50. [https://doi.org/10.1007/978-3-642-36809-7\\_2](https://doi.org/10.1007/978-3-642-36809-7_2).
45. Quinn TP, Erb I, Richardson MF, Crowley TM. Understanding sequencing data as compositions: an outlook and review. *Bioinformatics.* 2018;34(16):2870–8. <https://doi.org/10.1093/bioinformatics/bty175>.
46. Fernandes AD, Reid JN, Macklaim JM, McMurrugh TA, Edgell DR, Gloor GB. Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by compositional data analysis. *Microbiome.* 2014;2:15. <https://doi.org/10.1186/2049-2618-2-15>.
47. Mandal S, Van Treuren W, White RA, Eggesbø M, Knight R, Peddada SD. Analysis of composition of microbiomes: a novel method for studying microbial composition. *Microb Ecol Health Dis.* 2015;26: <https://doi.org/10.3402/mehd.v26.27663>.
48. Kuijjer ML, Tung MG, Yuan G, Quackenbush J, Glass K. Estimating Sample-Specific Regulatory Networks. *iScience.* 2019;14:226–40. <https://doi.org/10.1016/j.isci.2019.03.021>.
49. Nguyen T, Lee SC, Quinn TP, Truong B, Li X, Tran T, Venkatesh S, Le TD. Personalized Annotation-based Networks (PAN) for the Prediction of Breast Cancer Relapse. *bioRxiv.* 2019534628. <https://doi.org/10.1101/534628>.
50. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci.* 2005;102(43):15545–50. <https://doi.org/10.1073/pnas.0506580102>.
51. Xu L, Chen Y, Dutra-Clarke M, Mayakonda A, Hazawa M, Savinoff SE, Doan N, Said JW, Yong WH, Watkins A, Yang H, Ding L-W, Jiang Y-Y, Tyner JW, Ching J, Kovalik J-P, Madan V, Chan S-L, Müschen M, Breunig JJ, Lin D-C, Koeffler HP. *Proc Natl Acad Sci U S A.* 2017;114(15):3981–86. <https://doi.org/10.1073/pnas.1609758114>.
52. Sigismund S, Avanzato D, Lanzetti L. Emerging functions of the EGFR in cancer. *Mol Oncol.* 2018;12(1):3–20. <https://doi.org/10.1002/1878-0261.12155>.

## Publisher's Note

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