12-27-2022

# NetAct: a computational platform to construct core transcription factor regulatory networks using gene activity.

Kenong Su

Ataur Katebi

Vivek Kohar

Benjamin Clauss

Danya Gordin

*See next page for additional authors*

## Authors

Kenong Su, Ataur Katebi, Vivek Kohar, Benjamin Clauss, Danya Gordin, Zhaohui S Qin, Radha Krishna Murthy Karuturi, Sheng Li, and Mingyang Lu

## METHOD

# NetAct: a computational platform to construct core transcription factor regulatory networks using gene activity

Kenong Su[1†], Ataur Katebi[2,3†] , Vivek Kohar[4], Benjamin Clauss[3,5], Danya Gordin[2,3], Zhaohui S. Qin[6], R. Krishna M. Karuturi[7,8,9], Sheng Li[7,8] and Mingyang Lu[2,3,4,5*]

†Kenong Su and Ataur Katebi contributed equally to this work.

*Correspondence:
m.lu@northeastern.edu

[1] Department of Biomedical Informatics, Emory University, Atlanta, GA 30322, USA
[2] Department of Bioengineering|, Northeastern University, Boston, MA 02115, USA
[3] Center for Theoretical Biological Physics, Northeastern University, Boston, MA 02115, USA
[4] The Jackson Laboratory, Bar Harbor, ME 04609, USA
[5] Genetics Program, Graduate School of Biomedical Sciences, Tufts University, Boston, MA 02111, USA
[6] Department of Biostatistics and Bioinformatics, Emory University, Atlanta, GA 30322, USA
[7] The Jackson Laboratory for Genomic Medicine, Farmington, CT 06032, USA
[8] Department of Computer Science and Engineering, University of Connecticut, Storrs, CT, USA
[9] Graduate School of Biological Sciences & Eng., University of Maine, Orono, ME, USA

## Abstract

A major question in systems biology is how to identify the core gene regulatory circuit that governs the decision-making of a biological process. Here, we develop a computational platform, named NetAct, for constructing core transcription factor regulatory networks using both transcriptomics data and literature-based transcription factor-target databases. NetAct robustly infers regulators' activity using target expression, constructs networks based on transcriptional activity, and integrates mathematical modeling for validation. Our in silico benchmark test shows that NetAct outperforms existing algorithms in inferring transcriptional activity and gene networks. We illustrate the application of NetAct to model networks driving TGF-β-induced epithelial-mesenchymal transition and macrophage polarization.

**Keywords:** Systems biology, Gene regulatory networks, Gene regulatory circuits, Cellular state transitions, Mathematical modeling, Transcriptional activity, Epithelial-mesenchymal transition, Macrophage polarization

## Background

One of the major goals of systems biology is to infer and model complex gene regulatory networks (GRNs) which underlie the biological processes of human disease [1–6]. Particularly important are those gene networks that control decisions regarding cellular state transitions (e.g., replicative to quiescent [7–9], epithelial to mesenchymal (EMT) [10], pluripotent to differentiated [11, 12]), given the central importance of such regulatory processes to both healthy development as well as disease formation such as cancer tumorigenesis.

To construct and model GRNs associated with the biological process under investigation, researchers have developed two primary systems biology approaches. The first is a *bottom-up approach*, in which researchers focus on identifying a core GRN composed of a small set of master regulators [13]. Once the core GRN is obtained, mathematical

Su *et al. Genome Biology*    (2022) 23:270

Page 2 of 21

modeling is then applied to simulate the gene expression dynamics [14–17], which helps elucidate the potential gene regulatory mechanism driving the biological process in question. The current practice for synthesizing a core GRN is by compiling data via an extensive literature search, e.g., in these studies [18–20]. While this works well for systems where sufficient knowledge has been gained and accumulated, it is less effective in cases where key component genes and regulatory interactions have yet to be discovered. Due to the rapid increase of biomedical publications, manual synthesis of literature information has become extremely time-consuming and prone to human error in data interpretation. One way to address the labor-intensive issue is to rely on existing manually curated databases, such as KEGG [21] and Ingenuity Pathway Analysis (IPA) [22]. However, these databases often compile gene regulatory interactions from different tissues, species, or diseases. Therefore, it is hard to obtain context-specific interactions directly from these types of databases.

The second approach adopts a *top-down* perspective, in which researchers apply bioinformatics and statistical methods on genome-wide transcriptomics and/or genomics data to infer large-scale GRNs [13]. These data-driven methods are ideal for obtaining a global picture of gene regulation and the overall structure of gene-gene interactions. This approach can also be used to characterize key regulators and regulatory interactions between genes that are specific to the biological context of the study. However, conventional bioinformatics methods for gene network inference are usually not designed to identify an integrated working system. These methods typically rely on significance tests to determine the nodes and edges of a gene network, yet it is rare to evaluate whether the constructed gene network is capable of operating as a functional dynamical system [23]. Moreover, many statistical methods work well to identify the association between genes, but not their causation, thus limiting the applicative value of the top-down approach in characterizing gene regulatory mechanisms.

To overcome the abovementioned issues, a relatively new approach has been explored in several studies in which the top-down and bottom-up approaches are integrated to infer and model a core GRN [23–31]. In this combined approach, a GRN is constructed with bioinformatics tools using genome-wide gene expression data, followed by mathematical modeling of the GRN to simulate gene expression steady states and explore their similarity with biological cellular states. The simulations can help validate the accuracy of the constructed GRN and further clarify the regulatory roles of genes and interactions in driving cellular state transitions. This combined approach helps to discover existing and new regulatory interactions specific to the cell types and experimental conditions under study. Additionally, it helps pinpoint master regulators and reduce the system's overall complexity. The GRN modeling is particularly crucial for cases with non-trivial cellular state transitions, such as multi-step state transitions as observed in epithelial-mesenchymal transition (EMT) [32], and bifurcating state transitions, as observed in stem cell differentiation [33]. This is because the GRNs constructed by the top-down approach are not guaranteed to capture these state transition patterns. So far, to the best of our knowledge, there is no computational platform available that utilizes this combined approach for systematic GRN inference and modeling.

In this study, we introduce a computational platform, named NetAct, for inferring a core GRN of key transcription factors (TFs) using both transcriptomics data and a

Su *et al. Genome Biology*    (2022) 23:270

Page 3 of 21

literature-based TF-target database. Integrating both resources allows us to take full advantage of the existing knowledgebase of transcriptional regulation. NetAct adopts the combined top-down bioinformatics and bottom-up systems biology approaches, designed specifically to address the following two major issues.

First, many network inference methods rely on correlations of gene expression data, yet the actual transcriptional activities of many master regulators may not be reflected in their gene expression. Instead, the activity may be better associated with either their protein level, the level of a certain posttranslational modification, localization, or their DNA binding affinity. As a result, the master regulators with weak correlations between the expression level and the transcriptional activity will likely be discarded in the network. Some algorithms have been developed to infer the activities of regulators from transcriptomics data, such as VIPER [2], NCA [34], and AUCELL [35]. However, most of these algorithms (1) are not designed for gene network modeling, (2) still rely on the coexpression of a TF and its targeted genes, or (3) do not take advantage of the known regulatory interactions from the literature, hindering their applicability as automated algorithms for generic use in systems biology.
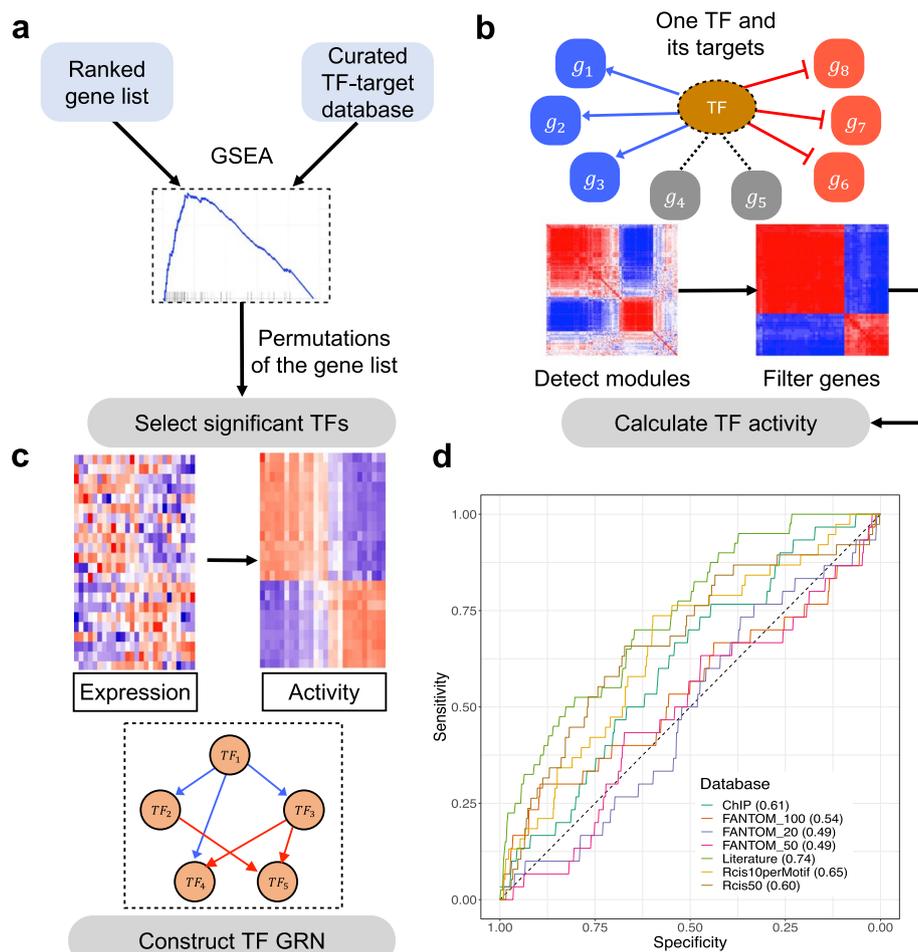
Second, conventional mathematical modeling approaches have been applied over the years to simulate the dynamics of a GRN, yet they are not particularly effective in analyzing core GRNs. A popular method models the gene expression dynamics of a system using the chemical rate equations that govern the associated gene regulatory processes. However, it is difficult to directly measure most of the kinetic parameters of a GRN. Although some parameter values can be learned from published results, many others are often based on educated guesses which significantly limits the predictive power of mathematical modeling. Moreover, a core GRN is not an isolated system. Thus, an ideal modeling paradigm should also consider other genes that interact with the core network. To address this infamous parameter issue, we have developed the modeling algorithm RACIPE [29, 36, 37] in previous work that analyzes a large ensemble of mathematical models with random kinetic parameters. RACIPE has been applied to model the dynamical behavior of gene regulatory networks of different biological processes, such as epithelial-mesenchymal transition [23, 29], cell cycle [37], and stem cell differentiation [38].

The new NetAct platform addresses the abovementioned issues by (1) inferring the activities of TFs for individual samples using the gene expression levels of their targeted genes, (2) identifying the regulatory interactions between two TFs based on their activities rather than their expressions, and (3) subsequently simulating the constructed core GRN with RACIPE to validate and evaluate the gene expression dynamics of the core GRN. In this paper, we describe in detail the NetAct platform, extensive benchmark tests for TF-target databases, TF activity inference, network construction, and two examples of applications to model GRNs with time series gene expression data.

## Results

We developed a computational systems biology platform, named NetAct, to construct transcription factor (TF)-based GRNs using TF activity. The method uniquely integrates both generic TF-target relationships from literature-based databases and context-specific gene expression data. NetAct also integrates our previously developed mathematical modeling algorithm RACIPE to evaluate whether the constructed network functions

Su *et al. Genome Biology*      (2022) 23:270

Page 4 of 21

properly as a dynamical system. It evaluates the roles of every gene in the network by in silico perturbation analysis. NetAct has three major steps: (1) identifying the core TFs using gene set enrichment analysis (GSEA) [39] with an optimized TF-target gene set database (Fig. 1a), (2) inferring TF activity (Fig. 1b), and (3) constructing a core TF network (Fig. 1c). Then, the network is validated and analyzed by simulating its dynamics using mathematical modeling by RACIPE. Details of each step are given in the "Methods" section and Additional file 1: Supplementary Note 5. Below, we demonstrate how we optimized the NetAct algorithm, compared its performance of activity inference with



**Fig. 1** Schematics of NetAct. **a** First, key transcription factors (TFs) are identified using gene set enrichment analysis (GSEA) with a literature-based TF-target database. **b** Second, the TF activity of an individual sample is inferred from the expression of target genes. From the co-expression and modularity analysis of target genes, we find target genes that are either activated (blue), inhibited (red), or not strongly related to the TF (gray). The activity is defined as the weighted average of target genes activated by the TF minus the weighted average of target genes inhibited by the TF. **c** Lastly, a TF regulatory network is constructed according to the mutual information of inferred TF activity and literature-based regulatory interactions. **d** Performance of GSEA for various TF-target gene set databases. The plot shows the sensitivity and specificity with different *q*-value cutoffs. The gene set databases in the benchmark include the combined literature-based database (D1); FANTOM5-based databases (D2) with 20, 50, and 100 target genes per TF; the combined experimental-based database (D3, ChIP); and RcisTarget databases (D4), one with 10 targets per TF binding motif and another with 50 total number of targets per TF

Su *et al. Genome Biology*      (2022) 23:270

Page 5 of 21

three existing methods using in silico gene expression data, and applied the network modeling approach to two biological datasets.

**Literature-based TF-target relationships facilitate TF inference**
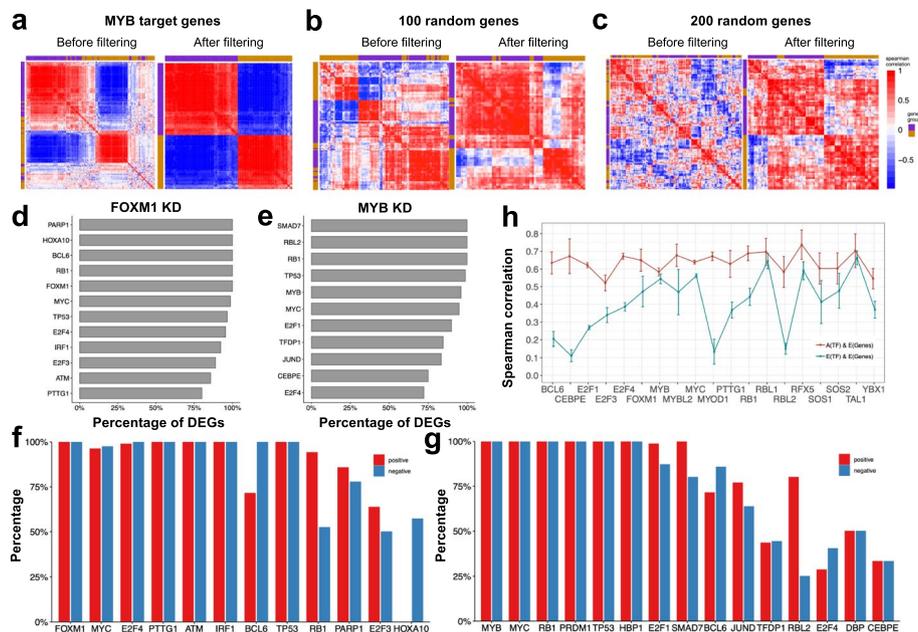
To establish a comprehensive gene set database containing TF-target relationships, we considered data from different sources (Additional file 3: Table S1, also see Additional file 1: Supplementary Note 1). They are (D1) a literature-based database, consisting of data from TRRUST [40], RegNetwork [41], TFactS [42], and TRED [43]; (D2) a gene regulatory network database FANTOM5 [44], whose interactions are extracted from networks constructed using RNA expression data from 394 individual tissues; (D3) a database derived from resources of putative TF binding targets, including ChEA [45], TRANSFAC [46], JASPAR [47], and ENCODE [48]; and (D4) a database derived from motif-enrichment analysis, RcisTarget [35]. These databases have been frequently used to study the transcriptional regulations and have already been utilized for network construction [29, 49].

We evaluated the performance of these databases by GSEA on a benchmark gene expression dataset. GSEA is a popular statistical method that can be used to evaluate the significant overlapping between a set of genes and differentially expressed genes between two experimental conditions. Using various types of TF-target databases, our goal is to find the best version of the database, so that GSEA can detect the target gene sets of the relevant TFs to be statistically significant. The benchmark dataset, denoted as *set B*, consists of a compilation of 12 microarray and 32 RNA-seq gene expression data (Additional file 3: Table S2). Each of these datasets contains at least three samples under the normal condition (control) and three samples under the treatment condition in which a specific TF is treated by knockdown (KD). We applied GSEA (with slight modifications, details in the "Methods" section) on set *B* to evaluate whether the enrichment analysis can detect the perturbed TFs. The underlying assumption is that, with a better TF-target gene set database, GSEA will be more likely to detect the corresponding perturbed TFs. For each TF-target database and each gene expression data in set *B*, we calculated the *q*-values of all the TFs in the database by GSEA to determine whether the target genes of the perturbed TF are enriched in the differentially expressed genes. We found that more significant *q*-values are usually associated with relatively larger number of targets for each TF; however, too many (e.g., greater than 2000) targets will result in non-significant *q*-values. The summary statistics, such as the total number of TFs and the average number of target genes per TF, are summarized in Additional file 3: Table S1. Furthermore, these corresponding *q*-values from all the gene expression data are converted to specificity and sensitivity values (see the "Methods" section), and different databases are compared based on the area under the sensitivity-specificity curves (Fig. 1d). We found that the literature-based database has the best overall performance; thus, we used this database for further analyses. Our results are in line with a previous benchmark study [50] that literature-based TF-target database outperforms others in capturing transcriptional regulation.

Su *et al. Genome Biology*      (2022) 23:270

Page 6 of 21

## Inferring TF activity without using TF expression

NetAct can accurately infer TF activity for an individual sample directly from the expression of genes targeted by the TF (see the "Methods" section). Here, we will illustrate how NetAct infers TF activity on two cases of microarray KD experiments—one case for shRNA KD of FOXM1 and shRNA KD of MYB in lymphoma cells (GEO: GSE17172 [51]), and another case for KD of BCL6 on both OCI-Ly7 and Pfeiffer GCB-DLBCL cell lines (GEO: GSE45838 [2]). NetAct first successfully identified the TFs that undergo knockdown in each case, i.e., FOXM1, MYB, and BCL6, by applying GSEA on the optimized TF-target database (*q*-value < 0.15).

Next, for each identified TF, NetAct calculates its activity using the mRNA expression of the direct targets of the TF. We first constructed a Spearman correlation matrix from the expression of the targeted genes. As shown in Fig. 2a, the correlation matrix after hierarchical clustering analysis typically consists of two red diagonal blocks, two blue off-diagonal blocks, and the remaining elements with low correlations which will be filtered out subsequently (details in the "Methods" section). Within the red blocks, the expression of any column gene is positively correlated with that of any row gene, while within the blue blocks, the expression of any column gene is negatively correlated with that of any row gene. This indicates that the genes in the two red blocks are



**Fig. 2** Illustration of the grouping scheme for target genes of a transcription factor. **a** The co-expression matrix of MYB target genes in shRNA knockdown of MYB lymphoma cells by hierarchical clustering analysis (Pearson correlation and complete linkage). **b**, **c** The poor clustering results from the co-expression of randomly selected 100 (**b**) and 200 genes (**c**). In panels **a**–**c**, the left subplots show the outcomes of all tested genes, and the right subplots show the outcomes of genes after the filtering step. Compared to the random cases, MYB target genes have a clear pattern of red and blue diagonal blocks from their co-expression. **d**, **e** The percentage of differentially expressed genes remained after the filtering step in the case of FOXM1 and MYB knockdown, respectively. **f**, **g** The proportion of genes from the activation group that are positively correlated with the TF expression (red bars) and the proportion of genes from the inhibition group that are negatively correlated with the TF expression (blue bars). **h** Spearman correlation (average and standard deviation) between TF activity and target expression (red) and between TF expression and target expression (blue)

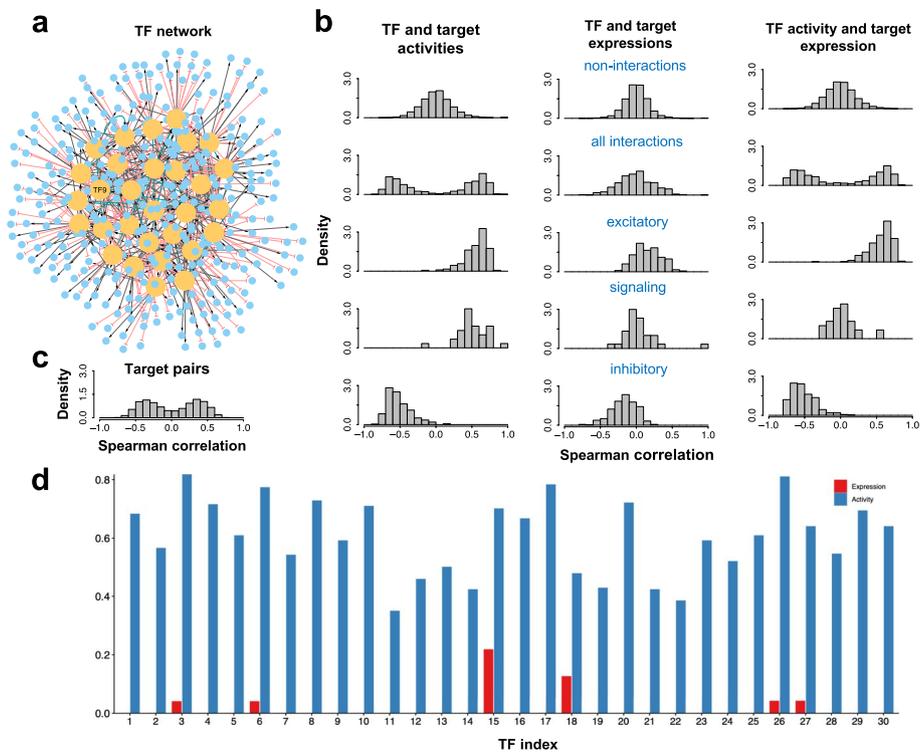Su *et al. Genome Biology*      (2022) 23:270

Page 7 of 21

anti-correlated in the gene expression with each other. However, if the correlation matrix is constructed from 100 or 200 randomly selected genes (Fig. 2b, c), such a clear pattern disappears. Thus, our observation suggests that genes from one of the red blocks are activated by the TF, whereas genes from the other block are inhibited by the TF. Moreover, filtered genes are not likely to be directly targeted by the TF in this context, or they are regulated by multiple factors simultaneously and are thus likely not a good indicator for the TF activity.

We further evaluated how the filtering step removes noise and retains the important genes in the analysis. We found that, after the filtering step, most of the differentially expressed (DE) genes are retained, as evidenced by Fig. 2d. Here, DE genes from each comparison were retrieved by using *limma* with a cutoff for the adjusted *p*-values at 0.05 and a cutoff for the log2 fold changes at 2. Subsequently, for DE TFs, we evaluated the Spearman correlations between the TFs and the corresponding targeted genes. In traditional approaches (such as ARACNe [1], WGCNA [52], and BEST [53]), the co-expression between a TF and its targeted genes is commonly used to identify its association and assign the sign (activation or inhibition) of the regulation. We found that, for each TF, most of the genes in a block either positively correlate with the TF expression (Fig. 2f, g, blue bars), or they negatively correlate with the TF expression (Fig. 2f, g, red bars). The tests demonstrate that, without directly using TF expression, NetAct can successfully identify two groups of important target genes—genes in each group are either activated or inhibited by the TF. These two groups of genes are further used to infer TF activity by a weighted average of their gene expression (Eq. 1 in the "Methods" section). Additionally, we found that the correlations between inferred TF activity and target expression are usually higher than the correlations between TF expression and target expression (Fig. 2h).

### Evaluating activity inference and network construction in a simulation benchmark

To evaluate the accuracy and robustness of inferred TF activity, we performed extensive benchmark tests to compare NetAct with other existing methods. We first performed the benchmark tests on simulated data because TF activity is usually not directly measurable. The activity of a TF can be related to its protein level or the level of a particular posttranslational modification, such as phosphorylation. Therefore, it is very difficult to obtain the ground truth of TF activity from an experimental dataset. Thus, in this benchmark test, we rely on mathematical modeling to simulate both the expression and activity of each TF from a synthetic TF-target network. With this simulated data, we benchmark NetAct against other methods.

To establish the simulated benchmark dataset, we first constructed a synthetic TF-target network with a total of 30 TFs. Each TF has 20 target genes randomly selected with replacement from a pool of 1000 genes. In addition, each TF also regulates two (randomly selected) of the 30 TFs. This synthetic network has a hierarchical structure, where a target gene may be co-regulated by multiple TFs. The type of each TF-to-TF regulation is either excitatory, inhibitory, or signaling, with a chance of 25%, 25%, and 50%, respectively; the type of each TF-to-target regulation is either excitatory or inhibitory with a 50% chance for each. Here, the signaling regulation changes the activity of a TF without changing its expression, whereas the excitatory or inhibitory interactions change both

Su *et al. Genome Biology*      (2022) 23:270

Page 8 of 21



**Fig. 3** Simulation of both gene expression and activity of a synthetic GRN. **a** The synthetic GRN consisting of 30 TFs and 447 target genes. An edge of transcriptional activation is shown as black line with an arrowhead; an edge of transcriptional inhibition as red line with a blunt head; an edge of signaling interaction as green line with an arrowhead. Transcription factor labeled as TF9 was selected for knockdown simulations. **b** The summary of the correlation analyses of the simulated expression and activity. The left, middle, and right columns represent the outcomes for TF and target activities, TF and target expressions, and TF activities and target expressions, respectively. For each category, the histograms of Spearman correlations are shown for non-interacting gene pairs (first row), interacting gene pairs (second row), gene pairs of excitatory transcriptional regulation (third row), gene pairs of excitatory signaling regulation (fourth row), and gene pairs of inhibitory transcriptional regulation (fifth row). Here, the target activity is set to be the same as the target expression for non-TF genes. **c** The histograms of Spearman correlations for gene pairs of target genes from the same TF. **d** Jaccard indices between the ground truth regulons of the synthetic GRN and the regulons inferred by ARACNe using either the simulated expression (red) or activity data (blue)

the activity and expression. From one realization of the synthetic network generation, the synthetic GRN contains a total of 477 genes (30 TFs, 447 targeted genes) and 660 regulatory links (Fig. 3a). See Additional file 1: Supplementary Note 4 for more details.
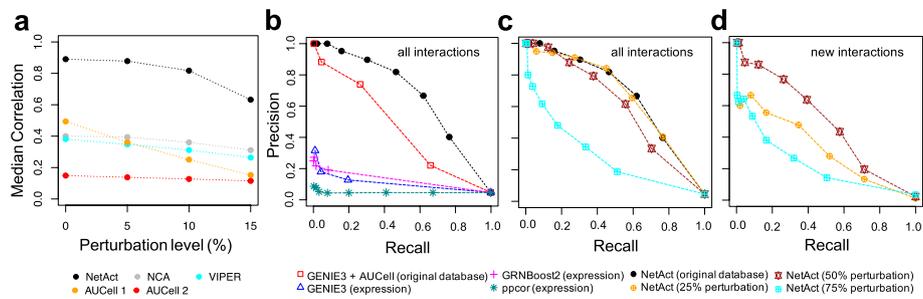
To simulate the gene expression of the TF-target network, we applied a generalized version of the mathematical modeling algorithm, RACIPE [37]. Using the network topology as the only input, RACIPE can generate an ensemble of random models, each corresponds to a set of randomly sampled parameters. Here, we used RACIPE to generate simulated data including gene expression and TF activity for the benchmark. Some previous studies have also adopted a similar modeling approach for benchmarking [54, 1]. To consider the effects of a signaling regulatory link, we generalized RACIPE to simulate both expression and activity for each TF. See Additional file 1: Supplementary Note 5 for more details.

In the benchmark test, we used RACIPE to simulate 100 models with randomly generated kinetic parameters. From these 100 models, we obtained 83 stable steady-state gene

Su *et al. Genome Biology*     (2022) 23:270

Page 9 of 21

expression and activity profiles for the 477 genes. As expected, TF activity and target activity from a regulatory link are correlated (1st column, 2nd row in Fig. 3b), TF activity and target expression (3rd column, 2nd row in Fig. 3b) are correlated, and the expression of two target genes (Fig. 3c) are correlated. However, there is no strong correlation between TF expression and target expression (2nd column, 2nd row in Fig. 3b) and, for a signaling regulatory link, between TF activity and target expression (3rd column, 4th row in Fig. 3b). Next, we applied ARACNe to predict the regulon (i.e., the list of targeted genes by a specific TF) using either the simulated expression profiles or the simulated activity profiles. We found that the regulons predicted from the activity profiles are substantially more similar to the predefined ground truth regulons (measured by the Jaccard index [55]) than those predicted from the expression profiles (Fig. 3d). The results indicate the need of using the TF activity, instead of TF expression, to identify TF-target relationships.

Next, we compared the performance of NetAct with several related algorithms, NCA, VIPER, and AUCell, in inferring TF activity using both the simulated expression profiles from the 83 models and a predefined regulon (i.e., the association of each TF with its target genes) (details for the implementation of these algorithms in Additional file 1: Supplementary Note 3). The predicted activity was then compared with the simulated activity (ground truth) to evaluate the performance. To mimic the real-life scenario where the target information may not be complete and accurate, we consider more challenging tests where the regulon data is randomly perturbed. Here, for a specific perturbation level, we generated 100 sets of regulon data by replacing a certain number of target genes for each TF with non-interacting genes. The numbers of replaced genes are 0 (0% level of perturbation), 5 (25%), 10 (50%), and 15 (75%) in different tests. We then evaluated the performance of NetAct, NCA, and VIPER. AUCell protocol advises to include the target genes with only positive interactions in the regulons. To satisfy this criterion, we updated the regulons for both unperturbed and perturbed regulons. For the unperturbed regulons, we retained only the positive interactions; for the perturbed regulons, we retained the positive target genes that were not replaced and a random half of the replaced target genes (assuming that half of the genes are positively regulated by the TF). We then evaluated AUCell performance using these updated regulons (denoted AUCell 1) and non-updated regulons (denoted AUCell 2). As shown in Fig. 4a (also Additional file 2: Figs. S3-S6), NetAct significantly outperforms each of the other methods in reproducing the simulated activity profiles at each perturbation level. As expected, the performance of NetAct is decreased by increasing the perturbation levels of the regulon data; however, NetAct still performs reasonably well even when only 25% of the actual target genes are kept in the regulon data. The results indicate that NetAct can robustly and accurately infer TF activity even with a noisy TF-target database.
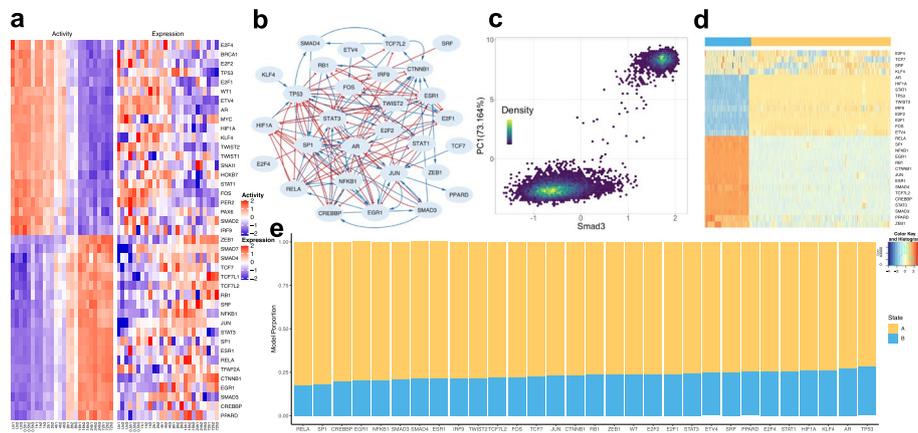
Furthermore, we tested another scenario where the test data contains simulated data from two experimental conditions, e.g., one representing an unperturbed condition and the other representing a perturbed condition. Here, we used the same synthetic network but compiled 40 expression and activity data from the abovementioned simulation (unperturbed condition), together with 43 expression and activity data from the simulations in which a specific TF (TF9) is knocked down (perturbed condition). We then performed a similar test as above and found that NetAct outperformed each of the other

Su *et al. Genome Biology*     (2022) 23:270

Page 10 of 21



**Fig. 4** The performance of activity and network inference from a simulation benchmark. **a** TF activity inference. TF activity was inferred by several methods using the gene expression data simulated from the synthetic TF-target gene regulatory network (GRN) and the corresponding regulons. For each TF, we computed Spearman correlations between the inferred activity and simulated activity (ground truth) for all the simulated models. Then, we calculated the average correlation values over all TFs. The plots show the median of average correlations for the cases where we used the original regulons defined by the TF-target network (0% perturbation), and the regulons where 5 (25% perturbation), 10 (50% perturbation), and 15 (75% perturbation) target genes are randomly replaced with non-interacting genes. The median values were computed over 100 repeats of random replacement for each perturbation level, and the values of the average correlations are reported for the case of zero perturbation. Shown are the results for NetAct (black), NCA (gray), VIPER (cyan), AUCELL 1 where regulons contain only positively associated target genes (orange), and AUCELL 2 where regulons contain all target genes (red). **b–d** Network inference. The panels show the performance of network inference algorithms from the simulation benchmark by the precision and recall for different link selection thresholds. **b** Network inference performance against all ground truth regulatory interactions. Tested methods are GENIE3, GRNBoost2, and PPCOR, using transcription factor (TF) expression; GENIE3 using TF activity inferred by AUCell; NetAct using its inferred TF activity. For the latter two methods, original (unperturbed) regulons obtained from the regulatory network were used. **c** Network inference performance of NetAct against all ground truth regulatory interactions using the regulons with 0% (the original), 25%, 50%, and 75% target perturbations. **d** Network inference performance of NetAct in discovering new regulatory interactions not existing in the regulons. NetAct was applied using the regulons at different perturbation levels (25%, 50%, and 75%). The benchmark results shown here are for the case of the untreated simulation. The results for the case of the knockdown simulation are shown in Additional file 2: Fig. S7

methods (Additional file 2: Fig. S2, Additional file 2: Fig. S7a). The notable performance gain of NetAct mainly emanates from the removal of incoherent (or noisy) targets of a TF before the activity calculation in NetAct (see the "Methods" section).

In addition, we performed a network construction benchmark of NetAct and a few other network construction algorithms using the in silico simulation dataset, as shown in Fig. 4b–d. NetAct, using the TF activity inferred from the original regulon database, outperforms not only network construction methods using gene expression, such as GENIE3 [56], GRNBoost2 [57], and ppcor [58, 59], but also GENIE3 using the TF activity inferred by AUCell (Fig. 4b). The last approach was presented to mimic a popular method SCENIC. Moreover, we evaluated the performance of NetAct when using a perturbed regulon database. We found that NetAct remains performing well when the perturbation level is as large as 50%, when evaluated by all the ground truth interactions (Fig. 4c) and by those not presented in the regulon database (Fig. 4d). The latter case was designed to evaluate the capability of NetAct in predicting novel interactions. We observed similar outcomes for the case of the second scenario of the simulation data from two conditions (Additional file 2: Fig. S7b-d, see Additional file 1: Supplementary Note 6 for details of the benchmark

Su *et al. Genome Biology*      (2022) 23:270

Page 11 of 21



**Fig. 5** Network modeling of TGF-β-induced EMT. Application of NetAct to an EMT in human cell lines using time-series microarray data. **a** Experimental expression and activity of enriched transcription factors. **b** Inferred TF regulatory network. Blue lines and arrowheads represent the gene activation; red lines and blunt heads represent gene inhibition. **c** The relationship between SMAD3 gene activity and the first principal component of the activity of all network genes from RACIPE simulations. **d** Hierarchical clustering analysis of simulated gene activity (with Pearson correlation as the distance function and Ward.D2 linkage method). Colors at the top indicate the two clusters from the simulated gene activity. The blue cluster represents the mesenchymal state, and the yellow cluster represents the epithelial state. The color legend for the heatmap is at the bottom right. **e** Knockdown simulations of the TF regulatory network. The bar plot shows the proportion of RACIPE models in each state (epithelial or mesenchymal) for the conditions of the knockdown of every TF

method). In summary, our in silico benchmark test demonstrates the high performance of NetAct over existing state-of-the-art methods in both inferring TF activity and gene regulatory networks.

**Characterizing cellular state transitions by GRN construction and modeling**

In the previous sections, we demonstrated the capability of NetAct in identifying the key TFs and predicting TF activity. With these data, NetAct further constructs a TF-based GRN using the mutual information (MI) of the activity from the identified TFs (details in the "Methods" section). We then applied RACIPE to the constructed network to check whether the simulated network dynamics are consistent with the experimental observations. Below, we show the utility of NetAct with two biological examples: epithelial-mechanical transition (EMT) and macrophage polarization.
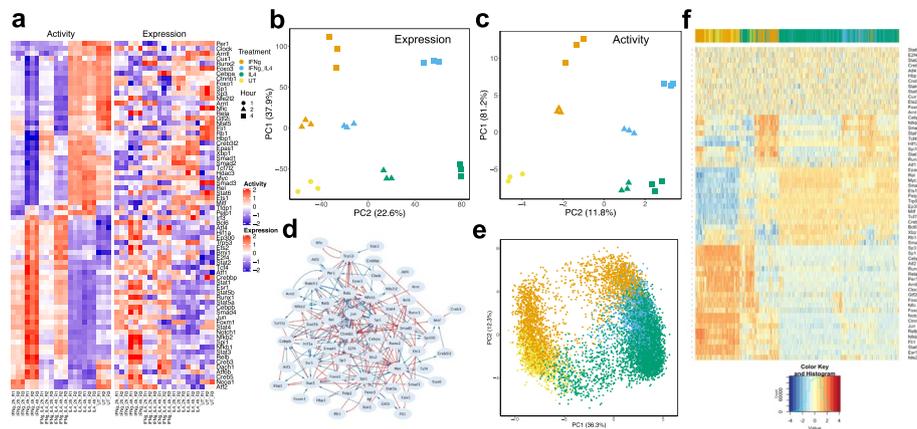
In the first case (EMT), we analyzed a set of time-series microarray data on A549 epithelial cells undergoing TGF-β-induced epithelial-mesenchymal transition (EMT) (GEO: GSE17708) [60]. According to the overall structure of the transcriptomics profiles, we arranged the samples from different time points into three groups—early stage (time points 0 h, 0.5 h, and 1 h), middle stage (time points 2 h, 4 h, and 8 h), and late stage (time points 16 h, 24 h, and 72 h). We then performed three-way GSEA with our human literature-based TF-target database to identify the enriched TFs that are active between early-middle, early-late, and middle-late time points. Forty-one TFs (*q*-value cutoff 0.01) were identified including many major transcriptional master regulators, such as BRCA1, CTNNB1, MYC, TWIST1, TWIST2, and ZEB1, and factors that are directly associated with TGF-β signaling pathways, such as SMAD3 [61], FOS, and JUN [62]. The hierarchical clustering analysis (HCA) of the expression and activity profiles for these TFs

is shown in Fig. 5a. While the expression profiles are quite noisy, the activities show a clear gradual transition from the epithelial (E) to mesenchymal (M) state. Note that the signs of the activity of a few non-DE TFs were flipped according to experimental evidence of protein-protein interactions and the nature of transcriptional regulation (see the "Methods" section for detailed procedures and Additional file 3: Table S3 for a list of the changes).

We then constructed a TF regulatory network (Fig. 5b) and performed mathematical modeling to simulate the dynamical behavior of the network using RACIPE (Fig. 5c, d). We found that, consistent with the expression and activity profiles (Fig. 5a), the network clearly allows two distinct transcriptional clusters that can be associated with E (the yellow cluster in Fig. 5d) and M states (the blue cluster in Fig. 5d). To assess the role of TGF-β signaling in inducing EMT, we performed a global bifurcation analysis [29] in which the SMAD3 level is used as the control parameter (Fig. 5c). Here, SMAD3 was selected as it is the direct target of TGF-β signaling [61]. As shown in Fig. 5c, when SMAD3 level is either very low or high, the cells reside in E or M states. However, when SMAD3 is at the intermediate level, the cells could be driven into some rare hybrid phenotypes. These results are consistent with our previous studies on the hybrid states of EMT [32, 63]. Using RACIPE, we systematically performed perturbation analyses by knocking down every TF in the network. Our simulation results (Fig. 5e) suggest that knocking down TFs, such as RELA, SP1, EGR1, and CREBBP, has major effects in driving M to E transition (MET), while knocking down TFs, such as TP53, AR, and KLF4, has major effects in driving E to M transition (EMT). These predictions are all consistent with existing experimental evidence (Additional file 3: Table S4).

Compared to a previous model of the EMT network based on an extensive literature survey [19], the GRN constructed by NetAct identified some of the same regulators induced by the TGF-β pathway, such as SMAD3/4, TWIST2, ZEB1, CTNNB1, NFKB1, RELA, FOS, and EGR1. Because of the lack of microRNAs and protein-protein interactions in the database, NetAct did not identify factors like miR200 and signaling molecules like PI3K. Interestingly, the NetAct model identifies STAT1/3, which was connected to other signaling pathways, such as HGF, PDGF, IGF1, and FGR, but not TGF-β in the previous network model. In addition, the NetAct model identified regulators in other important pathways in TGF-β-induced EMT in cancer cells, e.g., cell cycle pathway (RB1 and E2F1) and DNA damage pathway (P53).

In the second case, we studied the macrophage polarization program in mouse bone marrow-derived macrophage cells using time series RNA-seq data (GEO: GSE84517) [64]. In this experiment, macrophage progenitor cells (denoted as UT condition) were treated with (1) IFNγ to induce a transition to the M1 state, (2) IL4 to induce a transition to the M2 state, and (3) both IFNγ and IL4 to induce a transition to a hybrid M state. Here, we reprocessed the raw counts of RNA-seq with a standard protocol (details in Additional file 1: Supplementary Note 2). From principal component analysis (PCA) on the whole transcriptomics (Fig. 6b), we found that the gene expression undergoes distinct trajectories when macrophage cells were treated with either IFNγ (M1 state) or IL4 (M2 state). When both IFNγ and IL4 were administered, the gene expression trajectories are in the middle of the previous two trajectories, suggesting that cells are in a hybrid state (hybrid M state). We aim to use NetAct to elucidate the crosstalk in transcriptional

**Fig. 6** Network modeling of macrophage polarization. Application of NetAct to induced macrophage polarization via drug treatment in mice using RNA-seq data. **a** Experimental expression and activity of enriched TFs. **b** PCA projection of genome-wide gene expression profiles. Different point shapes indicate the time after treatment, and colors indicate treatment types **c** PCA projection of gene activity of enriched TFs. **d** Inferred TF regulatory network. Blue lines and arrowheads represent the gene activation; red lines and blunt heads represent the gene inhibition. **e** PCA projection of simulated gene activity of inferred network colored by mapping each model back to experimental data. **f** Hierarchical clustering analysis of simulated gene activity (with Pearson correlation as the distance function and Ward.D2 linkage method). Colors at the top indicate the mapped experimental conditions. The color legend of the heatmap is at the bottom

regulation downstream of cytokine-induced signaling pathways during macrophage polarization.

Here, we applied GSEA on six comparisons—untreated versus IFNγ-treated samples (one comparison between the untreated and the treated after 2 h, another between the untreated and the treated after 4 h, same for the other comparisons), untreated versus IL4-treated samples, and untreated versus IFNγ + IL4-treated samples. Using our mouse literature-based TF-target database, we identified 79 TFs (*q*-value cutoff 0.05 for UT vs IL4-2 h and 0.01 for all others). The expression and activity profiles of these TFs (Fig. 6a–c) capture the essential dynamics of transcriptional state transitions during macrophage polarization as follows. NetAct successfully identified important TFs in these processes, including Stat1, the major target of IFNγ, Stat2, Stat6, Cebpb, Nfkb family members, Hif1a, and Myc [65–67]. Myc is known to be induced by IL-4 at later phases of M2 activation and required for early phases of M1 activation [66]. Interestingly, we find Myc has high expression in both IL4 stimulation and its co-stimulation with IFN but its activity is high only in IL4 stimulation. We then constructed a TF regulatory network that connects 60 TFs (Fig. 6d) and simulated the network with RACIPE, from which we found that simulated gene expression (Fig. 6f) matches well with experimental gene expression data (Fig. 6a) (see Additional file 1: Supplementary Note 7). RACIPE simulations display disparate trajectories from UT to IL4 or IFNγ activation and stimulation with both IL4 and IFNγ. Strikingly, we found in the simulation that there is a spectrum of hybrid M states between M1 and M2 (Fig. 6e), which is consistent with the experimental observations of macrophage polarization [65]. Moreover, we also predict from our GRN modeling that the transition from UT to hybrid M is likely to first undergo a transition to either M1 or M2 before a second transition to hybrid M (Fig. 6e). This is because of our observation from the simulation data that there are fewer models connecting UT and hybrid M than any of the

other two routes (i.e., UT to M1, and UT to M2) (Additional file 2: Fig. S10). Taken together, we showed that the NetAct-constructed GRN model captures the multiple cellular state transitions during macrophage polarization.

In conclusion, we show that NetAct can identify the core TF-based GRN using both the literature-based TF-target database and the gene expression data. We also demonstrate how RACIPE-based mathematical modeling complements NetAct-based GRN inference in elucidating the dynamical behaviors of the inferred GRNs. Together, these two methods can be applied to infer biologically relevant regulatory interactions and the dynamical behavior of biological processes.

## Discussion

In this study, we have developed NetAct—a computational platform for constructing and modeling core transcription factor (TF)-based regulatory networks. NetAct takes a data-driven approach to establish gene regulatory network (GRN) models directly from transcriptomics data and takes a mathematical modeling approach to characterize cellular state transitions driven by the inferred GRN. The method specifically integrates both literature-based TF-target databases and transcriptomics data of multiple experimental conditions to accurately infer TF transcriptional activity based on the expression of their target genes. Using the inferred TF activity, NetAct further constructs a TF-based GRN, whose dynamics can then be evaluated and explored by mathematical modeling. Our approach in combining top-down and bottom-up systems biology approaches will contribute to a better understanding of the gene regulatory mechanism of cellular decision-making.

One of the key components of NetAct is a pre-compiled TF-target gene set database. Here, we have evaluated different types of TF-target databases in identifying knocked-down TFs using publicly available transcriptomics datasets. In this test, we have considered databases derived from the literature, gene co-expression, cis-motif prediction, and TF-binding motif data. Our benchmark tests suggest that the literature-based database clearly outperformed the other databases. The literature-based database usually contains a small (~ 30) number of target genes for each TF, but these data have direct experimental evidence, therefore being more reliable than those from the other sources. However, the literature-based database for sure has missing regulatory interactions, therefore maybe limiting the overall performance of NetAct. One way to address this issue is to further update the literature-based database, once new information is available. Another potential approach is to compile a database by combining different types of databases together. However, this might be quite challenging as different databases have data of very different sizes (the number of target genes) and quality. Future investigations on this direction can help to expand our knowledge of transcriptional regulation and meanwhile improve the performance of the algorithm.

NetAct also has a unique approach to infer the TF activity from the gene expression of the target genes with the consideration of activation/inhibition nature. From our in silico benchmark tests, we found that NetAct outperforms major activity inference methods, owing to the design of the filtering step and the use of a high-quality TF-target database. NetAct is also robust against some inaccuracy in the TF-target database and noises in

gene expression data, because of its capability of filtering out irrelevant targets as well as remaining key targets.

One potential issue is the assignment of the sign of TF activity, as it is algorithmically assigned according to the correlation with TF expression. In the case where the TF expression is very noisy or the expression is completely unrelated to TF activity, the sign assignment might be inaccurate. To deal with this issue, we have devised a semi-manual approach that identifies the sign of TF activity according to the sign of other interacting TFs. Another potential issue is that some TFs from the same family may have very similar target genes; therefore, NetAct will have difficulty in identifying exactly which TF from the family is most relevant. Additional data resources, such as epigenomics [68], TF-binding data [35], and Hi-C data [69], will be helpful to address this problem. One of the future directions is to design methods to integrate these data resources.

Lastly, instead of constructing a global transcriptional regulatory network, NetAct focuses on modeling a core regulatory network with only interactions between key TFs. The underlying hypothesis is that these TFs and the associated regulatory interactions play major roles in controlling the gene expression of different cellular states and the patterns of state transitions. With the core network identified using NetAct, we can further perform simulations with mathematical modeling algorithms, such as RACIPE, to analyze the control mechanism of the core network. These simulations allow us to generate new hypotheses, which can be further tested experimentally. The validation data can further help to improve the model. Ideally, this needs to be an iterative process to refine a core network model, which is indeed another interesting future direction.

## Conclusions

We developed NetAct, a computational platform for constructing and modeling core transcription factor regulatory networks using both transcriptomics data and literature-based transcription factor-target gene databases. Utilizing both types of resources allows us to identify regulatory genes and links specific to the data and fully take advantage of the existing knowledgebase of transcriptional regulation. Our method of combining top-down and bottom-up systems biology approaches contributes to a better understanding of the mechanism of gene regulation driving cellular state transitions.

## Methods

### Selecting enriched TFs

For a comparison between two experimental conditions, we obtained a ranked gene list quantified by the absolute value of the test statistics ($t$ statistics in microarray and Wald test statistics in RNA-seq) from differential expression (DE) analysis [70], followed by gene set enrichment analysis (GSEA) [39] using our optimized transcription factor (TF)-target gene set database. Here, for each TF, the corresponding gene set consists of all its target genes. GSEA identifies important TFs whose targets are enriched in DE genes between the two conditions. The significance test is achieved through 10,000 permutations of the gene list names and TFs are kept for further analysis when the $q$-value is below a certain threshold cutoff (0.05 by default). A C++ implementation of this version of GSEA, specifically for gene name permutations, has been provided in NetAct for fast

computation. For multiple comparisons, a set of enriched TFs are first identified from each pairwise comparison and then a union of the multiple sets of TFs is considered.

In the database benchmark test, for each database, we computed the sensitivity and specificity values for different $q$-value cutoffs. Here, for each cutoff value, we defined the sensitivity as the proportion of datasets where the gene sets for the KD TFs were enriched with $q$-values below the cutoff value. We also defined specificity as the fraction of cases where the gene sets for the other TFs (non-KD TFs in the benchmark) were not enriched with $q$-values above the cutoff value. We then computed the area under the ROC curve (AUC) using the DescTools R package [71].

### Inferring TF activity

TF activity is inferred from the expression of target genes retrieved from the TF-target database. NetAct defines the activity of the selected TFs using two different schemes—one using only the expression of target genes and the other using the expression of both the TF and its target genes. The second scheme is only used for the situation of noisy target gene expression. For each TF, the algorithm selects the better scheme according to its performance, as described below.

#### *Without directly using TF expression*

For each TF, its downstream targets are first divided into two modules using Newman's community detection algorithm [72] on the pairwise Spearman correlation matrix of the target genes. Then, within each module, some less-correlated genes are filtered out to improve the quality of the inference. Here, the filtering step is achieved as follows: (1) each target gene is assigned a vector of correlations with the other target genes, where the distance between two genes is calculated as the sum of squares of the correlation vectors of two genes; (2) $k$-mean algorithm ($k = 1$) is performed within each cluster to determine the center vector; and (3) genes are filtered out if the distance between the genes and the center is larger than the average distance.

This step outputs two groups of genes—genes in one group are supposed to be activated by the TF, while genes in the other group are inhibited by the TF. Note that, at this stage, the nature of activation/inhibition of the individual group is not yet determined. The activity of the TF is calculated as:

$$A(\text{TF}) = \frac{\sum_{i=1}^{n} w_i g_i I_i}{\sum_{i=1}^{n} w_i} \tag{1}$$

where $g_i$ is the standardized expression value of a target gene $i$, and $w_i$ is the weighting factor defined as a Hill function:

$$w_i = 1 / \left[ 1 + \left( \frac{s_i}{s_0} \right)^n \right] \tag{2}$$

where $s_i$ is the adjusted $p$-value from DE analysis for gene $i$, the threshold $S_0$ is 0.05, and $n$ is set to be 1/5 for best performance (Additional file 2: Fig. S8). $I_i$ is 1 if the corresponding gene belongs to the first group and $-1$ if it belongs to the second group. If the calculated TF activity pattern is not consistent with the TF expression trend (evaluated

by Spearman correlation), both the sign of the two groups and the sign of the activity are flipped. According to our in silico benchmark test (Additional file 2: Fig. S9), we found that majority of the targets in one group are activated by the TF, and majority of those in the other group are inhibited by the TF. For genes in the inhibition group, the higher the TF activity, the more the genes are suppressed. Thus, the formula in Eq. 1 captures well the activity of TFs for their effects to both activating and inhibitory targets. We also explored a few other community detection algorithms [73–75] and found they produced similar results (Additional file 2: Fig. S1).

### *Using TF expression*

For each TF, its downstream targets are first divided into two groups according to the sign of the Spearman correlation between the TF expression and the target expression. Similar to the previous scheme, in each group, target genes are filtered out if the correlation value is less than the average correlation of all the targets. The activity of the TF is also calculated using Eq. 1.

### *Sign assignment for DE TF*

For any DE TF (i.e., there is a significant difference in the TF expression across cell type conditions) of interest, NetAct computes the activity values from both the schemes (with or without TF's expression) and selects the better way based on how well the activity values correlate with target expression. To this end, NetAct calculates the absolute value of Spearman correlation between the TF activity and the expression of each target, and selects the scheme whose activity gives larger average correlations.

### *Sign assignment for non-DE TF*

If the expression patterns of the identified TFs fail to show the significant differences between cell type conditions, a semi-manual method to assign the sign of activity can be adopted. Putative interaction partners between DE and non-DE TFs in the inferred network are identified using Fisher's exact test between TF targets in the NetAct TF-target database. The most significant pairs are then cross-referenced with the STRING database (https://string-db.org) to identify instances of protein-protein interactions (PPIs). A literature search is then performed to identify the nature of the PPI, and the sign of the non-DE TF is adjusted based on the DE TF and the type of PPI. Note that the last step needs to be done manually for each modeling application. Additional file 3: Table S3 shows the details of TF sign flipping and supported experimental evidence for the two network modeling applications.

### Network construction and mathematical modeling

NetAct constructs a TF regulatory network using both the TF-TF regulatory interactions from the TF-target database and the activity values. (1) The network is constructed using mutual information between the activity values of two TFs. (2) Interactions are filtered out if they cannot be found in the TF-target regulatory database (i.e., D1). (3) The sign of each link is determined by the sign of the Spearman correlation between the activity of two TFs. (4) We keep the interaction between two

Su *et al. Genome Biology*     (2022) 23:270

Page 18 of 21

TFs if their mutual information is higher than a threshold cutoff. With different cut-off values for mutual information, NetAct establishes networks of different sizes. To identify the best network model capturing gene expression profiles, we apply mathematical modeling to each of the TF networks using RACIPE [29]. RACIPE takes network topology as the input and generates an ensemble of mathematical models with random kinetic parameters. By simulating the network, we expect to obtain multiple clusters of gene expression patterns that are constrained by the complex interactions in the network. RACIPE was also applied to generate simulated benchmark test sets for a synthetic TF-target network (Additional file 1: Supplementary Note 5).

## Supplementary Information

The online version contains supplementary material available at https://doi.org/10.1186/s13059-022-02835-3.

---

**Additional file 1: Supplementary Note 1.** TF-target gene set databases. **Supplementary Note 2.** Processing transcriptomics data. **Supplementary Note 3.** *In silico* TF activity inference benchmark. **Supplementary Note 4.** Construction of the synthetic GRN. **Supplementary Note 5.** Simulation of activity and expression using RACIPE. **Supplementary Note 6.** *In silico* network construction benchmark. **Supplementary Note 7.** Applications of network modeling with NetAct.

**Additional file 2: Fig. S1.** Comparison of the grouping schemes by Newman's method (NetAct) and other community detection algorithms. **Fig. S2.** Correlation structure in the simulated activities and expressions of the synthetic gene regulatory network with knockdown of transcription factor TF9. **Fig. S3.** Comparing stability of activity inference methods. **Fig. S4.** Null distribution of the average correlations for the four methods. **Fig. S5.** Activity levels of four transcription factors. **Fig. S6.** Scatter plot for activities of four transcription factors. **Fig. S7.** The performance of activity inference and network construction from a simulation benchmark. **Fig. S8.** Optimization of the Hill coefficient in the TF activity inference. **Fig. S9.** Comparison of NetAct grouping scheme for target genes with the synthetic gene regulatory network. **Fig. S10.** Analysis of 10,000 RACIPE-simulated gene expression profiles for the macrophage depolarization TF regulatory network.

**Additional file 3:** Supplementary tables: **Table S1.** Summary of the transcription factor (TF)-target databases for both human and mouse genomes. **Table S2.** Summary of the publicly available gene expression data sets for benchmarking TF-target gene set databases. **Table S3.** Sign correction for network construction and modeling. **Table S4.** Predicted driver TFs from network modeling of the EMT network and experimental evidences from the literature.

**Additional file 4:** Review history.

---

Su *et al. Genome Biology*        (2022) 23:270

Page 19 of 21

## 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.

## References

1.  Margolin AA, et al. ARACNE: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. BMC Bioinformatics. 2006;7:S7.
2.  Alvarez MJ, et al. Functional characterization of somatic mutations in cancer using network-based inference of protein activity. Nat Genet. 2016;48:838.
3.  Ament SA, et al. Transcriptional regulatory networks underlying gene expression changes in Huntington's disease. Mol Syst Biol. 2018;14:e7435.
4.  Chan TE, Stumpf MPH, Babtie AC. Gene regulatory network inference from single-cell data using multivariate information measures. Cell Syst. 2017;5:251–267.e3.
5.  Carré C, Mas A, Krouk G. Reverse engineering highlights potential principles of large gene regulatory network design and learning. Npj Syst Biol Appl. 2017;3:17.
6.  Fiers MWEJ, et al. Mapping gene regulatory networks from single-cell omics data. Brief Funct Genom. https://doi.org/10.1093/bfgp/elx046.
7.  Gérard C, Goldbeter A. Temporal self-organization of the cyclin/Cdk network driving the mammalian cell cycle. Proc Natl Acad Sci. 2009;106:21643–8.
8.  Laub MT, McAdams HH, Feldblyum T, Fraser CM, Shapiro L. Global analysis of the genetic network controlling a bacterial cell cycle. Science. 2000;290:2144–8.
9.  Li F, Long T, Lu Y, Ouyang Q, Tang C. The yeast cell-cycle network is robustly designed. Proc Natl Acad Sci. 2004;101:4781–6.
10. Nieto MA, Huang RY-J, Jackson RA, Thiery JP. EMT: 2016. Cell. 2016;166:21–45.
11. Kim J, Chu J, Shen X, Wang J, Orkin SH. An extended transcriptional network for pluripotency of embryonic stem cells. Cell. 2008;132:1049–61.
12. Loh Y-H, et al. The Oct4 and Nanog transcription network regulates pluripotency in mouse embryonic stem cells. Nat Genet. 2006;38:431.
13. Katebi A, Ramirez D, Lu M. Computational systems-biology approaches for modeling gene networks driving epithelial–mesenchymal transitions. Comput Syst Oncol. 2021;1:e1021.
14. Alon U. An introduction to systems biology: design principles of biological circuits. (Chapman and Hall/CRC); 2006. https://doi.org/10.1201/9781420011432.
15. Kirk PDW, Babtie AC, Stumpf MPH. Systems biology (un)certainties. Science. 2015;350:386–8.
16. Chasman D, Roy S. Inference of cell type specific regulatory networks on mammalian lineages. Curr Opin Syst Biol. 2017;2:130–9.
17. Ben-Jacob E, Lu M, Schultz D, Onuchic JN. The physics of bacterial decision making. Front Cell Infect Microbiol. 2014;4:154.
18. Dutta P, Ma L, Ali Y, Sloot PMA, Zheng J. Boolean network modeling of β-cell apoptosis and insulin resistance in type 2 diabetes mellitus. BMC Syst Biol. 2019;13:36.
19. Steinway SN, et al. Network modeling of TGFβ signaling in hepatocellular carcinoma epithelial-to-mesenchymal transition reveals joint Sonic Hedgehog and Wnt pathway activation. Cancer Res. 2014;74:5963–77.
20. Zeigler AC, et al. Computational model predicts paracrine and intracellular drivers of fibroblast phenotype after myocardial infarction. Matrix Biol. 2020;91–92:136–51.
21. Kanehisa M, Furumichi M, Sato Y, Ishiguro-Watanabe M, Tanabe M. KEGG: integrating viruses and cellular organisms. Nucleic Acids Res. 2021;49:D545–51.
22. Krämer A, Green J, Pollard J, Tugendreich S. Causal analysis approaches in Ingenuity Pathway Analysis. Bioinforma Oxf Engl. 2014;30:523–30.
23. Ramirez D, Kohar V, Lu M. Toward modeling context-specific EMT regulatory networks using temporal single cell RNA-Seq data. Front Mol Biosci. 2020;7:54.
24. Dunn S, Li MA, Carbognin E, Smith A, Martello G. A common molecular logic determines embryonic stem cell self-renewal and reprogramming. EMBO J. 2019;38:e100003.
25. Wooten DJ, Gebru M, Wang H-G, Albert R. Data-driven math model of FLT3-ITD acute myeloid leukemia reveals potential therapeutic targets. J Pers Med. 2021;11:193.
26. Udyavar AR, et al. Novel hybrid phenotype revealed in small cell lung cancer by a transcription factor network model that can explain tumor heterogeneity. Cancer Res. 2017;77:1063–74.
27. Wooten DJ, et al. Systems-level network modeling of small cell lung cancer subtypes identifies master regulators and destabilizers. PLoS Comput Biol. 2019;15:e1007343.
28. Khan FM, et al. Unraveling a tumor type-specific regulatory core underlying E2F1-mediated epithelial-mesenchymal transition to predict receptor protein signatures. Nat Commun. 2017;8:198.

Su *et al. Genome Biology*      (2022) 23:270

Page 20 of 21

29. Kohar V, Lu M. Role of noise and parametric variation in the dynamics of gene regulatory circuits. Npj Syst Biol Appl. 2018;4:1–11.
30. Moignard V, et al. Decoding the regulatory network of early blood development from single-cell gene expression measurements. Nat Biotechnol. 2015;33:269–76.
31. Sha Y, Wang S, Zhou P, Nie Q. Inference and multiscale model of epithelial-to-mesenchymal transition via single-cell transcriptomic data. Nucleic Acids Res. 2020;48:9505–20.
32. Lu M, Jolly MK, Levine H, Onuchic JN, Ben-Jacob E. MicroRNA-based regulation of epithelial–hybrid–mesenchymal fate determination. Proc Natl Acad Sci. 2013;110:18144–9.
33. Jang S, et al. Dynamics of embryonic stem cell differentiation inferred from single-cell transcriptomics show a series of transitions through discrete cell states. eLife. 2017;6:e20487.
34. Liao JC, et al. Network component analysis: reconstruction of regulatory signals in biological systems. Proc Natl Acad Sci. 2003;100:15522–7.
35. Aibar S, et al. SCENIC: single-cell regulatory network inference and clustering. Nat Methods. 2017;14:1083–6.
36. Huang B, et al. Interrogating the topological robustness of gene regulatory circuits by randomization. PLoS Comput Biol. 2017;13:e1005456.
37. Katebi A, Kohar V, Lu M. Random parametric perturbations of gene regulatory circuit uncover state transitions in cell cycle. iScience. 2020;23:101150.
38. Huang B, et al. Decoding the mechanisms underlying cell-fate decision-making during stem cell differentiation by random circuit perturbation. J R Soc Interface. 2020;17:20200500.
39. Subramanian A, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. Proc Natl Acad Sci U S A. 2005;102:15545–50.
40. Han H, et al. TRRUST: a reference database of human transcriptional regulatory interactions. Sci Rep. 2015;5:11432.
41. Liu Z-P, Wu C, Miao H, Wu H. RegNetwork: an integrated database of transcriptional and post-transcriptional regulatory networks in human and mouse. Database. 2015;2015.
42. Essaghir A, Demoulin J-B. A minimal connected network of transcription factors regulated in human tumors and its application to the quest for universal cancer biomarkers. PLoS One. 2012;7:e39666.
43. Jiang C, Xuan Z, Zhao F, Zhang MQ. TRED: a transcriptional regulatory element database, new entries and other development. Nucleic Acids Res. 2007;35:D137–40.
44. Abugessaisa I, et al. FANTOM5 transcriptome catalog of cellular states based on Semantic MediaWiki. Database J Biol Databases Curation. 2016;2016.
45. Lachmann A, et al. ChEA: transcription factor regulation inferred from integrating genome-wide ChIP-X experiments. Bioinformatics. 2010;26:2438–44.
46. Wingender E, Dietze P, Karas H, Knüppel R. TRANSFAC: a database on transcription factors and their DNA binding sites. Nucleic Acids Res. 1996;24:238–41.
47. Sandelin A, Alkema W, Engström P, Wasserman WW, Lenhard B. JASPAR: an open-access database for eukaryotic transcription factor binding profiles. Nucleic Acids Res. 2004;32:D91–4.
48. Luo Y, et al. New developments on the Encyclopedia of DNA Elements (ENCODE) data portal. Nucleic Acids Res. 2020;48:D882–9.
49. Abugessaisa, I. et al. FANTOM5 transcriptome catalog of cellular states based on Semantic MediaWiki. Database J. Biol. Databases Curation 2016, baw105 (2016).
50. Garcia-Alonso L, Ibrahim MM, Turei D, Saez-Rodriguez J. Benchmark and integration of resources for the estimation of human transcription factor activities.  Genome Res. 2021;31(4):745.
51. Alvarez MJ, Sumazin P, Rajbhandari P, Califano A. Correlating measurements across samples improves accuracy of large-scale expression profile experiments. Genome Biol. 2009;10:R143.
52. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. BMC Bioinformatics. 2008;9:559.
53. Hu M, Qin ZS. Query large scale microarray compendium datasets using a model-based Bayesian approach with variable selection. PLoS One. 2009;4:e4495.
54. Schaffter T, Marbach D, Floreano D. GeneNetWeaver: in silico benchmark generation and performance profiling of network inference methods. Bioinformatics. 2011;27:2263–70.
55. Levandowsky M, Winter D. Distance between sets. Nature. 1971;234:34.
56. Huynh-Thu VA, Irrthum A, Wehenkel L, Geurts P. Inferring regulatory networks from expression data using tree-based methods. PLoS One. 2010;5:e12776.
57. Moerman T, et al. GRNBoost2 and Arboreto: efficient and scalable inference of gene regulatory networks. Bioinformatics. 2019;35:2159–61.
58. Kim S. ppcor: an R package for a fast calculation to semi-partial correlation coefficients. Commun Stat Appl Methods. 2015;22:665–74.
59. Pratapa A, Jalihal AP, Law JN, Bharadwaj A, Murali TM. Benchmarking algorithms for gene regulatory network inference from single-cell transcriptomic data. Nat Methods. 2020;17:147–54.
60. Sartor MA, et al. ConceptGen: a gene set enrichment and gene set relation mapping tool. Bioinformatics. 2010;26:456–63.
61. Schiffer M, Von Gersdorff G, Bitzer M, Susztak K, Böttinger EP. Smad proteins and transforming growth factor-β signaling. Kidney Int. 2000;58:S45–52.
62. Zhang Y, Feng X-H, Derynck R. Smad3 and Smad4 cooperate with c-Jun/c-Fos to mediate TGF-β-induced transcription. Nature. 1998;394:909–13.
63. Jolly MK, et al. Implications of the hybrid epithelial/mesenchymal phenotype in metastasis. Front Oncol. 2015;5.
64. Piccolo V, et al. Opposing macrophage polarization programs show extensive epigenomic and transcriptional crosstalk. Nat Immunol. 2017;18:530–40.
65. Mosser DM, Edwards JP. Exploring the full spectrum of macrophage activation. Nat Rev Immunol. 2008;8:958–69.
66. Bae S, et al. MYC-mediated early glycolysis negatively regulates proinflammatory responses by controlling IRF4 in inflammatory macrophages. Cell Rep. 2021;35:109264.

Su *et al. Genome Biology*     (2022) 23:270

Page 21 of 21

67. Hu X, Ivashkiv LB. Cross-regulation of signaling pathways by interferon-γ: implications for immune responses and autoimmune diseases. Immunity. 2009;31:539–50.
68. Pliner HA, et al. Cicero predicts cis-regulatory DNA interactions from single-cell chromatin accessibility data. Mol Cell. 2018;71:858–871.e8.
69. Malysheva V, Mendoza-Parra MA, Saleem M-AM, Gronemeyer H. Reconstruction of gene regulatory networks reveals chromatin remodelers and key transcription factors in tumorigenesis. Genome Med. 2016;8:57.
70. Ritchie ME, et al. limma powers differential expression analyses for RNA-sequencing and microarray studies. Nucleic Acids Res. 2015;43:e47.
71. Signorell, A. et al. DescTools: tools for descriptive statistics. (2022).
72. Newman MEJ. Modularity and community structure in networks. Proc Natl Acad Sci. 2006;103:8577–82.
73. Reichardt J, Bornholdt S. Statistical mechanics of community detection. Phys Rev E. 2006;74:016110.
74. Newman MEJ. Analysis of weighted networks. Phys Rev E. 2004;70:056131.
75. Newman MEJ. Finding community structure in networks using the eigenvectors of matrices. Phys Rev E. 2006;74:036104.
76. Su K, Katebi A, Kohar V, Clauss B, Gordin D, Qin Z, et al. NetAct analysis code and data. GitHub. 2022. https://github.com/lusystemsbio/NetActAnalysis.
77. Su K, Katebi A, Kohar V, Clauss B, Gordin D, Qin Z, et al. NetAct analysis code and data. GitHub (Zenodo link). 2022. https://doi.org/10.5281/zenodo.7352281.
78. Su K, Katebi A, Kohar V, Clauss B, Gordin D, Qin Z, et al. NetAct R package. GitHub. 2022; https://github.com/lusystemsbio/NetAct.
79. Su K, Katebi A, Kohar V, Clauss B, Gordin D, Qin Z, et al. NetAct R package GitHub (Zenodo link); 2022. https://doi.org/10.5281/zenodo.7352299.

## Publisher's Note