

Dynamic regulatory module networks for inference of cell type specific transcriptional networks

Alireza Fotuhi Siahpirani^{1,2,+}, Deborah Chasman^{1,8,+}, Morten Seirup^{3,4}, Sara Knaack¹, Rupa
Sridharan^{1,5}, Ron Stewart³, James Thomson^{3,5,6}, and Sushmita Roy^{1,2,7*}

¹Wisconsin Institute for Discovery, University of Wisconsin-Madison

²Department of Computer Sciences, University of Wisconsin-Madison

³Morgridge Institute for Research

⁴Molecular and Environmental Toxicology Program, University of Wisconsin-Madison

⁵Department of Cell and Regenerative Biology, University of Wisconsin-Madison

⁶Department of Molecular, Cellular, & Developmental Biology, University of California Santa
Barbara

⁷Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison

⁸Present address: Division of Reproductive Sciences, Department of Obstetrics and Gynecology,
University of Wisconsin-Madison

+These authors contributed equally.

*To whom correspondence should be addressed.

Abstract

Changes in transcriptional regulatory networks can significantly alter cell fate. To gain insight into transcriptional dynamics, several studies have profiled transcriptomes and epigenomes at different stages of a developmental process. However, integrating these data across multiple cell types to infer cell type specific regulatory networks is a major challenge because of the small sample size for each time point. We present a novel approach, Dynamic Regulatory Module Networks (DRMNs), to model regulatory network dynamics on a cell lineage. DRMNs represent a cell type specific network by a set of expression modules and associated regulatory programs, and probabilistically model the transitions between cell types. DRMNs learn a cell type's regulatory network from input expression and epigenomic profiles using multi-task learning to exploit cell type relatedness. We applied DRMNs to study regulatory network dynamics in two different developmental dynamic processes including cellular reprogramming and liver dedifferentiation. For both systems, DRMN predicted relevant regulators driving the major patterns of expression in each time point as well as regulators for transitioning gene sets that change their expression over time.

Introduction

Transcriptional regulatory networks connect regulators such as transcription factors to target genes, and specify the context specific patterns of gene expression. Changes in regulatory networks can significantly alter the type or function of a cell, which can affect both normal and disease processes. The regulatory interaction between a transcription factor (TF) and a target gene's promoter is dependent upon TF binding activity, histone modifications and open chromatin, that have all been associated with differential expression between cell types [1–4]. To probe the dynamic and cell type-specific nature of mammalian regulatory networks, several research groups are generating matched transcriptomic and epigenomic data from short time courses or for cell types related by a branching lineage [5–7]. However, very few methods have exploited these datasets to infer cell type specific regulatory networks.

Existing computational methods to infer cell type specific networks can be grouped into three main categories: (i) Skeleton network-based methods, (ii) Regression-based methods, (iii) Probabilistic graphical model-based methods. *Skeleton-network based methods*, combine available protein-protein and protein-DNA interactions to create a “skeleton” network representing the union of possible edges that can exist in a cell at any time, and overlay context-specific mRNA levels on this network to derive dynamic snapshots of the network [8, 9]. Such approaches rely on a comprehensive characterization of the regulatory network, which we currently lack for most mammalian cell types. A second group uses *linear and non-linear regression-based methods* to predict mRNA levels as a function of chromatin marks [10, 11] and/or transcription factor occupancies [11] and can infer a predictive model of mRNA for a single condition (time point or cell type). These regression approaches are applied to each context individually and have not been extended to model multiple related time points or cell types, which is important to study how networks transition between different time points and cell states. *Multi-task regression methods* [12–17] are systematic approaches to learn multiple related networks. Thus far, these approaches have nearly all been based on mRNA levels and require a sufficiently large number of mRNA samples for each time point or cell type to reliably estimate the statistical dependency structure. The last class of methods are based probabilistic graphical models, namely, dynamic Bayesian networks (DBNs), including input-output Hidden Markov Models [18] and time-varying DBNs [15]. Both DREM and time-varying DBNs are suited for time courses only, and do not accommodate lineage trees.

To address the limitations of existing methods and systematically integrate parallel transcriptomic and epigenomic datasets, we have developed a novel dynamic network reconstruction method, Dynamic Regulatory Module Networks (DRMNs) to predict regulatory networks in a cell type-specific manner by leveraging their relationships to each other. DRMNs are based on a non-stationary probabilistic graphical model and can be used to model regulatory networks on a lineage. DRMNs represent the cell type regulatory network by a concise set of gene expression modules, defined by groups of genes with similar expression levels, and their associated regulatory programs. The module-based representation of regulatory networks enables us to reduce the number of parameters to be learned and increases the number of samples available for parameter estimation. To learn the regulatory programs of each module at each time point, we use multi-task learning that shares information between related time points or cell types.

We applied DRMN to three datasets measuring transcriptomic and epigenomic profiles in multiple cell types at different stages of cellular reprogramming and during dedifferentiation from liver hepatocytes. Two of these corresponding to cellular reprogramming including one array [19] and one sequencing experiments [7]. The third dataset (unpublished) includes RNA-seq and ATAC-seq profiling for mouse dedifferentiation. DRMN learned a modular regulatory program for each of the cell types by integrating chromatin marks, open chromatin, sequence specific motifs and gene expression. Compared to an approach that does not model dependencies among cell types DRMN was able to predict more realistic regulatory networks that were more concise and reflected the shared relatedness among the cell types, while maintaining high predictive power of expression. Furthermore, integrating cell type specific chromatin data with cell type invariant sequence motif data enabled us to better predict expression than each data type alone. Comparison of the inferred regulatory networks showed that they change gradually over time, and identified key regulators that are different between the cell types. In particular, among the regulators identified by DRMN were Klf4, Myc, and Tcf3 which are known to be associated with ESC state. In the hepatocyte data, we found regulators associated with development such as Gbx2, Six6, and Irx1 associated with transitioning gene sets. Taken together our results show that DRMN is a powerful approach to infer cell type specific regulatory networks, which enables us to systematically link upstream regulatory programs to gene expression states and to identify regulator and module transitions associated with changes in cell state.

Results

Dynamic Regulatory Module Networks (DRMNs)

DRMNs are used to represent and learn cell-type or time-point specific regulatory networks while leveraging the relationship among the cell types and time points. DRMNs design is motivated by the difficulty of inferring a cell type specific regulatory network from expression when only a small number of samples are available for each cell type. In a DRMN, the regulatory network for each context, such as a cell type, is represented compactly by a set of modules, each representing a discrete expression state, and the regulatory program for each module (**Figure 1**). The regulatory program is a predictive model of expression that predicts the expression of genes in a module from upstream regulatory features such as sequence motifs and epigenomic signals. Like previous module-based representations of regulatory networks [14,20,21], DRMN organizes genes into modules to reduce the number of parameters and increase the number of samples available for statistical network inference. In addition, DRMN leverages the relationship between cell types defined by a lineage tree. The modules and regulatory programs are learned simultaneously using a multi-task learning framework to encourage similarity between the regulators chosen for a cell type and its parent in the lineage tree. DRMN allows two ways to share information across cell types: DRMN-FUSED and DRMN-ST. DRMN-FUSED makes use of regularized regression to share information across cell types, while the other uses a graph structure prior approach (See **Methods**). Both approaches are able to share information across timepoints effectively, however, the DRMN-FUSED approach is computationally more efficient.

DRMNs offer a flexible framework to integrate diverse regulatory genomic features

DRMN's predictive model can be used to incorporate different types of regulatory features, such as sequence-based motif strength, accessibility, and histone marks. We first examined the relative importance of context-specific (e.g., chromatin marks and accessibility) and context independent features (e.g., sequence motifs) for building an accurate gene expression model. We compared the performance of both versions of DRMNs on different feature sets on a mouse reprogramming dataset which had both array and sequencing datasets (**Figure 2**). Our metric for comparison was the Pearson's correlation between true and predicted expression

in each module (**Figure 2A**) on three fold cross-validation.

The different features we considered were, (1) Motif only, (2) Histone marks, (3) Histone and Motif, (4) Accessibility (from atac-seq) and motif, (5) Motifs with quantified accessibility (Q-motif), (6) Histone marks, motif features and accessibility, (7) Histone marks, Q-motif, and (8) Histone marks, Q-motif and accessibility. We first compared motif alone, chromatin alone, and chromatin+motif, as these features were available for both array and sequencing datasets for the DRMN-ST (**Figure 2B,C**) and DRMN-FUSED (**Figure 2D,E**) models. In both models, motifs alone (dark blue markers) have minimal predictive power, which is consistent across different k and for both array (**Figure 2B, D**) and sequencing (**Figure 2C, E**) data. Histone marks alone (red marker) have higher predictive power, however adding both histone marks and motif features (magenta) has the best performance for $k = 3$ and 5, with the improved performance to be more striking for DRMN-ST. For DRMN-Fused, Histone only and Histone + motifs seemed to perform similarly, although at higher k using histone marks is better. Between different cell types the performance was very consistent in array data, while for sequencing data, the MEF and MEF48 cell types were harder to predict than ESC and preIPSC.

We next examined the contribution of accessibility data in predicting expression. As only the sequencing dataset included ATAC-seq data for each cell type, we performed this comparison on the sequencing dataset alone (**Figure 2C, E**). We incorporated the ATAC-seq data in five ways: as a single feature defined by the aggregated accessibility of a particular promoter (ATAC feature, orange markers **Figure 2**) combined with motif sequence, using ATAC-seq to quantify the strength of a motif instance (Q-Motif, cyan markers **Figure 2B, D**), combining the ATAC feature with histone and motifs (dark purple marker), combining Q-Motif with histone (light green), and the ATAC feature with histone and Q-motifs (dark green). We also considered ATAC-seq as a single feature, but this was not very helpful (**Supplementary Figure S1**).

Combining the ATAC feature together with motif feature improves performance over the sequence motif feature (**Figure 2B, D** orange *vs.* dark blue markers) suggesting that these are both useful features for predicting expression. The Q-Motif feature (cyan markers) was better than the motif only feature (dark blue) at lower k ($k=3$), however, surprisingly, did not outperform the sequence alone features. One possible explanation for this is that the Q-Motif feature was more sparse than the sequence only feature, due to the 0 feature value if we could not assign a read to the feature. Finally, we compared the ATAC feature combined

with chromatin and motif features to study additional gain in performance (**Figure 2B, D**, purple markers). Interestingly, even though using ATAC-seq features combined with motif features improves the performance over motif alone, addition of ATAC-seq feature to histone marks+motif does not change the performance of either version of DRMN (**Figure 2B, D**) compared to histone + motif (magenta markers). This suggests that combination of multiple chromatin marks capture the dynamics of expression levels better than the chromatin accessibility signal. It is possible that the overall cell type specific information captured by the accessibility profile is redundant with the large number of chromatin marks and we might observe a greater benefit of ATAC-seq if there were fewer or no marks. We observe similar trends with Histone+Q-Motif and Histone+ATAC+Q-motif features (light and dark green) which perform on par to each other, and close to histone+motif and histone+ATAC+motif.

DRMN-ST and DRMN-FUSED behaved in a largely consistent way for these different feature combinations with the exception of the Histone+Motif feature where DRMN-Fused was not gaining in performance at higher k . However, when we directly compared DRMN-ST to DRMN-FUSED, DRMN-FUSED was able to outperform DRMN-ST on most of the feature combinations (**Figure 2F, G**), with the exception of Histone+Motif (magenta), Histone + Motif + ATAC (dark purple), Histone + Q-Motif (light green), and Histone + Q-Motif + ATAC (dark green) where the two models were similar. It is likely that DRMN-ST learns a sparser model at the cost of predictive power (**Supplementary Figure S2**). For the application of DRMNs to real data, we focus on DRMN-FUSED due to its improved performance.

Multi-task learning approach is beneficial for learning cell type-specific expression patterns

We next compared the utility of DRMN to share information across cell types or time points while learning predictive models of expression against several baseline models: (1) those that do not incorporate sharing (RMN), or (2) those that are clustering-based (GMM-Merged and GMM-Indep). The clustering-based baselines offer simple approaches to describe the major expression patterns across time but link *cis*-regulatory elements to expression changes only as post-processing steps. For both DRMN and RMN, we learned predictive models of expression in each module using different regulatory feature sets (**Figure 3, Supplementary Figure S3**): motif only, chromatin mark only and using motif and chromatin marks. We used overall correlation and per module expression level comparison to assess the quality of the predictions. We

performed these experiments on the array and the RNA-seq dataset for mouse reprogramming.

Using the overall correlation, both variants of DRMN, RMN model and GMM-Indep, vastly outperform GMM-Merged, which also had the least stable results across cell types and folds (**Supplementary Figure S3A-F**). This suggests that the gene partitions are likely different between the different cell types and imposing a single structure for all three, as done in GMM-Merged, misses out on the cell type specific aspects of the data. Based on overall correlation, DRMN models performed at par with RMN and GMM-Indep for most cases (**Supplementary Figure S3A-F**), with the exception of Histone and Histone+Motif for sequencing data (**Supplementary Figure S3A, B**) where GMM-Indep is worse for lower k s. These results suggest that based on overall correlation, clustering could offer a first approach to analyze these data, however, a predictive modeling approach has advantages over a clustering approach as it improves prediction quality by incorporating additional *cis*-regulatory data.

We next compared the different models on the basis of the per-module expression levels (**Supplementary Figure S3G-L**). DRMN and RMN clearly outperform the GMM-based approaches, which is expected because GMM is only able to produce one value per module and does not capture the within-module variation. The overall high correlation for DRMN and RMN demonstrate that a model learned from the regulatory features is able to provide a more fine-tuned model of expression variation.

Finally, we compared DRMNs to RMNs (**Figure 3**). Both versions of DRMN outperform or are at par with their corresponding RMN versions on histone only and histone + motif features (**Figure 3**, blue for DRMN and red for RMN). On motifs, the difference between the models was dependent upon the cell line, the number of modules, k and specific implementation of DRMNs. In particular, DRMN-Fused was at par or better than RMs on sequence motifs. However, DRMN-ST had a lower performance for several cell lines on the sequencing data, at higher k . Interestingly, DRMN-ST did have a greater gain in performance compared to DRMN-Fused on the Histone-Motif and Histone datasets from arrays. It is likely that motif only features are being overfit to the data resulting in reduced performance. Overall, our results suggest that a predictive modeling approach as in DRMN and RMN, can explain the expression variation much better than a clustering based approach, and that using multi-task learning with regularized regression helps to build accurate models of gene expression.

Using DRMN to gain insight into regulatory programs of cellular reprogramming

We applied DRMN to gain insights into the regulatory programs of cellular reprogramming. We focus on the results obtained on the sequencing dataset for reprogramming (**Figure 4**), although many of the trends are captured in the array data too (**Supplementary Figure S4**). DRMN modules learned on the sequencing data exhibit seven distinct patterns of expression in each of the four cell types (**Figure 4A**). We observe that while the expression patterns remains the same, the number of genes in each module in each cell type varies (**Figure 4A**). Furthermore, we compared the extent of similarity of matched modules between consecutive cell types and found that the modules were on average 30-90% similar, exhibiting the lowest similarity at the MEF48 to pips transition and the highest between the MEFs (**Figure 4C**). For the array as well, we observed the greatest dissimilarity between the MEF and pips transition (**Supplementary Figure S4**). This agrees with the pips state exhibiting a major change in transcriptional status during reprogramming.

To interpret the modules, we next looked at the regulatory programs inferred for each module (**Figure 4B**). While some of these regulators are selected in all cell types, we also observe some cell type specific patterns, such as *Pitx2* in module 4 (associated with muscle cell differentiation [22]), *Pou5f1* in module 5 and *Esrrb* in module 6 (associated with pluripotency state [3]), and *Insm1* in module 7 (associated with differentiation [23], all highlighted in red in **Figure 4B**).

To gain insight into the dynamics of the process, we identified genes that change their module assignments between time points. Overall, of the 17,358 genes, 11,104 were associated with a module transition in the sequencing dataset (**Figure 5A**) and of the 15,982 total genes in the array dataset, 3,573 exhibited a module transition (**Supplementary Figure S4A**). We clustered the transitioning genes and tested each for biological processes and also *cis*-regulatory elements. These transitioning gene sets can provide insight into the overall dynamics of the process. One of these sets is enriched for binding sites of Tcf3 (**Figure 5B**), which is known to repress pluripotency genes [24, 25] and is only enriched in the low expression module (module 1) in MEF and MEF48. This gene set, induced specifically in ESCs was predicted to be regulated by histone marks H3K79me2 and H3K27me3 and the TF, Bcl6, which indicates an interplay of histone and TF binding to enable cell fate specification. These genes move from low expression module (1) in early stages (MEF and MEF 48) to highly expressed modules (5, and 6) in the ES state (cell type specific module assignments on left, cell type specific expression in the middle). We also observe that the promoter regions

of the genes in this set are strongly enriched for the Tcf3 motif (dark purple heatmap on right). The genes exhibiting this trend include several ESC specific genes such as *Sall4* and *Dppa2*. We observe a similar pattern with the array data as well (**Supplementary Figure S5**), identifying ESC specific genes induced in the iPSC state and enriched for Tcf3, thus corroborating the transitioning gene sets between the two expression platforms.

Using DRMNs to gain insight into regulatory programs of hepatocyte dedifferentiation

We next applied DRMNs to understand the temporal dynamics of regulatory programs during hepatocyte dedifferentiation. A major challenge in studying dynamics in primary cells such as liver hepatocytes is that they dedifferentiate from their hepatocyte state. Maintaining hepatocytes in their differentiated state is important for studying normal liver function as well as for liver-related diseases [26]. Dedifferentiation could be due to the changes in the regulatory program over time, however, little is known about the transcriptional and epigenetic changes during this process. To address this, we collected RNA-seq and ATAC-seq data for 16 time points from 0 hours to 36 hours (**Figure 6**).

Here we applied DRMNs with $k=5$ modules and followed a similar approach for first interpreting the modules by regulatory programs inferred for each module. We observe pluripotency associated TFs such as *Gbx2* and *Kl4* [27], TFs associated with hepatocyte differentiation such as *Onecut2* [28], and liver associated TFs such as *Hnf4a* [29]. We also tested these genes for GO enrichment (**Supplementary Figure S6**) and found that the repressed module (Module 1) was enriched for developmental processes while the other modules were enriched for diverse metabolic processes. Of these modules 4 and 5, which are associated with higher expression are enriched for more liver-specific metabolic function such as co-enzyme metabolism, acetyl CoA metabolism and modules 2 and 3 where enriched for general housekeeping function such as DNA and nucleic acid metabolism.

We next examined the transitioning gene sets. We identified a total of 150 transitioning gene sets spanning 5,762 genes. Many of the transitions were between modules that are adjacent to each other based on expression levels, suggesting the majority of the dynamic transitions are subtle (e.g., module 1 and 2, **Figure 7A**). We next predicted regulators for these transitioning gene sets using a regularized regression model, MTG-LASSO (**Methods**). Using this approach we identified 84 gene sets that we could predict a regulator

of. **Figure 7B** shows 5 of these gene sets where module assignment change more than 1 (e.g. 1 to 2 to 3) and regulators associated with these gene sets. These include pluripotency or development associated TFs such as Irx1 [30] (gene set 420 and 439), Gbx2 [27, 31] and Six6 [32] (gene set 439), Mecom [33], Esrrb [3], Alx1 [34] (gene set 390). Together, these results suggest that analysis regulatory interactions associated with specific modules and transitioning gene sets, can help identify transcription factors that drive these processes.

Discussion

Cell type-specific gene expression patterns are the result of a complex interplay between transcription factor binding, genome accessibility and histone modifications. Accordingly, datasets that measure both transcriptomes and epigenomes in closely related cell types and processes are becoming increasingly available. A key challenge is to interrogate these datasets to identify the underlying gene regulatory networks that drive context-specific expression changes. However, doing so from such datasets is a major challenge because of the large number of variables measured in each time point. In this work, we developed, DRMNs, that simplifies genome scale regulatory networks into gene modules and infers regulatory program for each module in all the input cell types. Using DRMNs, one can characterize the major transcriptional patterns during a dynamic process over time or over a cell lineage and identify transcription factors and epigenomic signals that are responsible for these transitions.

Central to DRMN's modeling framework is to jointly learn the regulatory programs for each cell type or time point by using multi-task learning. Using two different approaches to multi-task learning, we show that joint learning of regulatory programs is advantageous compared to a simpler approach of learning regulatory programs independently per cell type. Furthermore, predictive modeling of expression that also clusters genes into expression groups is more powerful than learning a single predictive model and simple clustering. Such models have improved generalizability and are able to capture fine-grained expression variation as a function of the upstream regulatory state of a gene.

DRMN offers a flexible framework to integrate a variety of experimental setups. In its simplest form, DRMN can be applied to an expression time course, array or sequencing, and can use sequence specific motif instances to learn a predictive model. However, DRMN is able to integrate other types of regulatory signals such as chromatin accessibility measured using ATAC-seq, histone modifications and transcription factor measured using ChIP-seq. The datasets that we applied DRMN to include exemplars of different designs. In particular the reprogramming sequencing dataset was the most comprehensive with nine histone marks, ATAC-seq and expression measured for four cell types. In contrast the dedifferentiation dataset measured RNA-seq and ATAC-seq for 16 time points.

The Chronis et al [7] dataset enabled us to systematically study the utility of different cell-type specific measurements such as chromatin marks and accessibility to predict expression. When combining ATAC-seq

with chromatin marks to predict expression, we did not see a substantial improvement in predictive power. This is likely because of the large number of chromatin marks in our dataset. However, in both array and sequencing data we found that combine sequence features (motifs) and histone modifications had the highest predictive power. In our application to the Chronis et al dataset, we did not observe a substantial advantage of using accessible motif instances (Q-motif) as opposed to all motif instances. This is most likely due to the sparsity of the feature space of Q-motifs. As future work, an important direction of work would be to explore additional ways to incorporate the ATAC-seq signal for each motif could provide a greater benefit [2].

We applied DRMNs to two distinct types of developmental processes: mouse reprogramming (3-4 cell types) and hepatocyte dedifferentiation (16 time points). DRMN application identified the major patterns of expression in both sequencing and array reprogramming datasets as well as the dedifferentiation dataset. When comparing the results from both processes there were greater changes in expression in the reprogramming time course compared to dedifferentiation indicative of the different dynamics in the two processes. Importantly DRMN was able to identify key regulators for each module as well as dynamic gene sets in both cases. Among the regulators identified by DRMN for the induced genes in the reprogramming dataset included known pluripotency markers. Similarly, DRMN predicted liver transcription factors such as HNF4G::HNF4A for induced genes in the dedifferentiation dataset. In both datasets, we identified transitioning genes and predicted regulators for these genes using simple or multi-task regression. This enabled us to identify regulators important for transitions. In particular for ESC, we found a gene set that exhibited dynamics predictable by histone elongation marks and a transcription factor Bcl6. Similarly for hepatocyte dedifferentiation, we found gene sets associated with changes in HNF4G motif accessibility.

Currently DRMN operates on regulatory features on gene promoter regions. An important direction of future work would be to incorporate long-range interactions to enable distal regions to contribute to the expression levels of gene. Another direction would be to use more generic sequence features, such as k-mers [35] to enable the discovery of novel regulatory elements and offer great flexibility in capturing sequence specificity and its role in predictive models of expression.

Taken together, DRMN offers a powerful and flexible framework to model time series and lineage specific regulatory genomic datasets and enables inference of cell type-specific regulatory programs. As datasets that profile epigenetic and transcriptomic dynamics of specific processes become available, methods

like DRMNs will become increasingly useful to examine cell type and context-specific regulatory networks form these datasets.

Materials and methods

Dynamic Regulatory Module Networks (DRMNs)

DRMN model description

The DRMN is a probabilistic model of regulatory networks for multiple cell types related by a lineage tree. DRMN is suited for datasets with multiple time points or conditions, with a small number of samples (e.g, one or two) per condition and several types of measurements, such as RNA-seq, ChIP-seq and ATAC-seq. Due to the small sample size, a standard regulatory network, that connects individual TFs to target genes for each time point, cannot be inferred. Instead in DRMN, we learn regulatory programs for groups of genes. For C cell types, the DRMN model is defined by a set of regulatory module networks, $\mathbf{R} = \{R_1, \dots, R_C\}$, a lineage tree τ , and transition probability distributions $\mathbf{\Pi} = \{\Pi_1, \dots, \Pi_C\}$, **Figure 1**. For each cell type c , $R_c = \langle G_c, \Theta_c \rangle$ defines the regulatory program as G_c , the set of edges between regulators and modules, and Θ_c , the parameters of a regression function for each module that relates the selected regulatory features to the expression of the module's genes. Transition matrices $\{\Pi_1, \dots, \Pi_C\}$ capture the dynamics of the module assignments in one cell type c given its parent cell type in the lineage tree τ . Specifically, $\Pi_c(i, j)$ is the probability of any gene being in module j given that its parental assignment is to module i . For the root cell type, this is simply a prior probability over modules.

The DRMN inputs are (i) cell type-specific expression for N genes $\mathbf{X} \in \mathcal{R}^{N \times C}$, assuming we have a single measurement of a gene in each cell type; (ii) gene-specific regulatory features for F candidate regulators $\mathbf{Y}^{N \times C \times F}$; and (iii) the lineage tree, τ . The regulatory signals used in the regulatory program can be either context-independent (e.g., a sequence-based motif network) or context-specific (e.g., a motif network informed by epigenomic measurements). The number of modules, k , is provided as input to the method. Given the above inputs, the DRMN model aims to optimize the following score:

$$P(\mathbf{R}|\mathbf{X}, \mathbf{Y}) \propto P(\mathbf{X}|\mathbf{R}, \mathbf{Y})P(\mathbf{R}) \quad (1)$$

Here $P(\mathbf{X}|\mathbf{R}, \mathbf{Y})$ is the data likelihood given the regulatory program, and $P(\mathbf{R})$, is a prior probability distribution of the regulatory program. The data likelihood can be decomposed over the individual cell types: $P(\mathbf{X}|\mathbf{R}, \mathbf{Y}) = \prod_c P(X_c|R_c, Y_c)$. Within each cell type, this is modeled using a mixture of predictive models, one model for each module. The prior, $P(\mathbf{R})$ can in turn be defined in terms of the graph structure G_c and parameters, Θ_c . We used two formulations for this prior to enable sharing information: DRMN-Structure Prior (DRMN-ST) defines a structure prior over the graph structures $P(G_1, \dots, G_C)$ while DRMN-FUSED uses a regularized regression framework and implicitly defines priors on the $P(\theta_1, \dots, \theta_C)$. In both frameworks, we share information between the cell types/time points to learn the regulatory programs of each cell type or time point.

DRMN-ST: Structure prior approach. The structure prior, $P(\mathbf{R})$, which is defined only over the graph structures, $P(G_1, \dots, G_C)$ assumes the parameters are set to their Maximum likelihood setting. The prior term $P(G_1, \dots, G_C)$ determines how information is shared between different cell types at the level of the network structure and encourages similarity of regulators between cell types. $P(G_1, \dots, G_C)$ is computed using the transition matrices $\{\Pi_1, \dots, \Pi_C\}$ and decomposes over individual regulator-module edges within each cell type as follows:

$$P(G_1, \dots, G_C) = \prod_{f \rightarrow k} P(\mathbf{I}_{f \rightarrow k})$$

where

$$P(\mathbf{I}_{f \rightarrow k}) = P(I_{f \rightarrow k}^{root}) \prod_{c' \rightarrow c \in \tau} P(I_{f \rightarrow k}^c | I_{f \rightarrow k}^{c'})$$

where $I_{f \rightarrow k}^c$ is an indicator function for the presence of the edge $f \rightarrow k$ for cell type c , between a regulator f and a module k . To define $P(I_{f \rightarrow k}^c | I_{f \rightarrow k}^{c'})$, we use the transition probability of the modules as:

$$P(I_{f \rightarrow k}^c | I_{f \rightarrow k}^{c'}) = \begin{cases} \Pi_c(k|k) & \text{if } I_{f \rightarrow k}^c = I_{f \rightarrow k}^{c'} \\ 1 - \Pi_c(k|k) & \text{otherwise} \end{cases}$$

Here the first option gives the probability of maintaining the same state from parent to child cell lines, if the edge has the same state in both parent and child (is present in both c and c' or is absent in both), and the

second option gives the probability of changing the edge state.

We implemented a greedy hill climbing algorithm to incorporate the structure prior. See Section **DRMN learning** for more details of the learning algorithm.

DRMN-FUSED: Parameter prior approach. Our second approach used a fused group LASSO formulation to share information between the cell types/time points by defining the following objective:

$$\min_{\Theta} \sum_c \|X_{c,k} - Y_{c,k} \Theta_{c,k}^T\|_2^2 + \rho_1 \|\Theta_k\|_1 + \rho_2 \|R\Theta_k\|_1 + \rho_3 \|\Theta_k\|_{2,1}, \quad (2)$$

where $X_{c,k}$ is the expression vector of genes in module k in cell line c , $Y_{c,k}$ is the feature matrix (size of module k by F , the number of features) corresponding to the same genes, and $\Theta_{c,k}$ is the vector of regression coefficients for the same module and cell line (1 by F , non-zero values correspond to selected features), and Θ_k is the C by F matrix resulting from concatenation of $\Theta_{c,k}$ vectors (number of cell types/time points by number of features, C by F). R is a $C - 1$ by C matrix, encoding the lineage tree. Each row of R correspond to an edge of the tree, meaning that if row i of R correspond to the edge $c_1 \rightarrow c_2$ in the lineage tree, we set $R(i, c_1)$ to 1, $R(i, c_2)$ to -1 , and all other values in that row to 0. This would mean that $R\Theta_k$ would be a $C - 1$ by F where row i correspond to $\Theta_{c_1,k} - \Theta_{c_2,k}$, the difference between regression coefficients of cell lines c_1 and c_2 . $\|\cdot\|_1$ denotes l_1 -norm (sum of absolute values), $\|\cdot\|_2$ denotes l_2 -norm (square root of sum of square of value), and $\|\cdot\|_{2,1}$ denotes $l_{2,1}$ -norm (sum of l_2 -norm of columns of the given matrix). ρ_1, ρ_2 and ρ_3 correspond to hyper parameters, with ρ_1 for sparsity penalty, ρ_2 to enforce similarity between selected features of consecutive cell lines in the lineage tree, and ρ_3 to enforce selecting the same features for all cell lines/time points. Thus, ρ_2 controls the extent to which more closely related cell types are closer in their regulatory program, which ρ_3 controls the extent to which all the cell types share similarity in their regulatory programs. We set these hyper parameters based on cross-validation by performing a grid search for $\rho_1 \in \{0.5, 1, 2, 5, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, 120, 130\}$, and $\rho_2, \rho_3 \in \{0, 10, 20, 30, 40, 50\}$. We implemented the algorithm described in MALSAR Matlab package [36], which uses accelerated gradient method [37, 38] to minimize the objective function above in C++.

DRMN learning

DRMNs are learned by optimizing the DRMN score (**Eqn 1**), using an Expectation Maximization (EM) style algorithm that searches over the space of possible graphs for a local optimum (**Algorithm 1**). The algorithm uses a multi-task learning approach to jointly learn the regulatory programs for all cell types. In the M step, we estimate transition parameters (M1 step) and the regulatory program structure (M2 step). In the E step, we compute the expected probability of a gene's expression profile to be generated by one of the regulatory programs.

Algorithm 1: DRMN Algorithm

Input:

- Expression data $\mathbf{X} = \{X_1, \dots, X_C\}$,
- Regulatory features $\mathbf{Y} = \{Y_1, \dots, Y_C\}$,
- Initial module assignments $\mathbf{M} = \{M_1, \dots, M_C\}$

Output:

- Regulatory programs $\mathbf{R} = \{R_1 = (G_1, \Theta_1), \dots, R_C = (G_C, \Theta_C)\}$,
- Transition probabilities $\mathbf{\Pi} = \{\Pi_1, \dots, \Pi_C\}$

while not converged do

M1: Estimate transition parameters (Π_1, \dots, Π_C)

M2: Update regulatory programs ($G_1, \dots, G_C, \Theta_1, \dots, \Theta_C$)

E: Update module assignment probabilities and module assignments ($\mathbf{\Gamma}, M_1, \dots, M_C$)

end

Estimate transition parameters (M1 step): Let $\gamma_{k,k'}^{g,c}$ be the probability of gene g in cell type c to belong to module k , and, in its parent cell type c' , to module k' . We calculate the probability of transitioning from k' in c' to k in c as $\Pi_c(k, k') = \frac{\sum_g \gamma_{k,k'}^{g,c}}{\sum_{g,k,k'} \gamma_{k,k'}^{g,c}}$.

Update regulatory programs (M2 step): Recall that the regulatory program for each cell type c is $R_c = \langle G_c, \Theta_c \rangle$, where G_c is a set of regulatory interactions $f \rightarrow k$ from a regulatory feature f to a module k , and Θ_c are the parameters of a regression function for each module that relates the selected regulatory features to the expression of the module's genes. The estimation of the regulatory programs is specific to the way in which information is shared across tasks, and differs in the DRMN-ST and DRMN-FUSED approach.

We assume that the expression levels are generated by a mixture of experts, each expert corresponding to a module. Each expert uses a Multivariate normal distribution, which means that the likelihood of gene

g 's expression being generated by the regulatory program for a module k in cell type c can be written as:

$$\begin{aligned} P(x_{g,c}|R_{c,k}, \mathbf{y}_{g,c}) &\sim N(\mu_{c,k}, \Sigma_{c,k}) \\ &= N(\theta_{c,k}^0 + \sum_{f \in G_{c,k}} \theta_{c,k}^f y_{g,c}^f, \Sigma_{c,k}) \end{aligned}$$

where $G_{c,k}$ is the current list of regulators selected for module k in cell type c . This is the equivalent of a linear regression to predict expression of genes in module k in cell type c , using regulatory features of genes in cell type c . In the DRMN-ST approach, the regulatory interactions are learned by per-module per-cell type regression in a greedy hill-climbing framework. At initialization, all regulatory network structures G_c are empty, and the parameters Θ_c are computed simply as the empirical mean and variance of the genes initially assigned to each module in each cell type. In each iteration of DRMN learning, we fix the current module assignments M_c for all cell types c and update the regulatory program of each module independently. In each iteration, we score each potential regulatory feature based on its improvement to the likelihood of the model, and choose the regulator with maximum improvement (if over a minimum threshold). This regulator is added to the module's regulatory program for all cell types for which it improves the cell type-specific likelihood. For the regularized regression version, DRMN-FUSED, the parameters are learned using an accelerated gradient method [37, 38] to minimize the objective function in **Eqn 2**.

Update soft module assignments (E step): Let $\gamma_{k|k'}^{g,c}$ be the probability of gene g in cell type c to belong to module k , given that in its parent cell type c' , it belonged to module k' . We also introduce α_g^c , a vector of size $K \times 1$ where each element $\alpha_g^c(c')$ specifies the probability of observations given the parent state is c' . We estimate the probabilities using a dynamic programming procedure, where values at internal nodes in the lineage tree are computed using the values for all descendent nodes, down to the leaves.

If c is a leaf node, we calculate

$$\gamma_{k|k'}^{g,c} = P(x_{g,c}|R_{c,k}, \mathbf{y}_{g,c}) \Pi_c(k|k')$$

where the first term is the probability of observing expression of gene g in cell type c in module k (given its

regulatory program and regulatory features), and the second is the probability of transitioning from module k' in the parent cell type c' to module k in cell type c .

For a non-leaf cell type c :

$$\gamma_{k|k'}^{g,c} = P(x_{g,c}|R_{c,k}, \mathbf{y}_{g,c})\Pi_c(k|k') \prod_{c \rightarrow l \in \tau} \alpha^{g,l}(k)$$

For both internal and leaf cell types, we write the joint probability of the (k', k) pair as $\gamma_{k,k'}^{g,c} = \frac{\gamma_{k|k'}^{g,c}}{\alpha^{g,c}(k')}$, where $\alpha^{g,c}(k') = \sum_k \gamma_{k|k'}^{g,c}$ is the probability of g 's expression in any module given parent module k' .

Termination: DRMN inference runs for a set number of iterations or until convergence. Final module assignments are computed as maximum likelihood assignments using a dynamic programming approach. While module assignments between consecutive iterations do not change significantly, the final module assignments are significantly different from the initial module assignments, and predictive power of model significantly improves as iterations progress (though improvements are small after 10 iterations, **Supplementary Figure S7**).

In our experiments, we ran DRMN for up to ten iterations. When using greedy hill climbing approach, to decrease computation time, we inferred regulatory edges in small batches per DRMN iteration. For each full DRMN iteration, the M2 step was run until up to five regulators were added per module.

Hyper parameter tuning of DRMN-FUSED To assess the effect of hyper-parameters on the performance of DRMN when using fused group LASSO, we performed a grid search on a wide range of parameters values. **Supplementary Figure S8** shows the average correlation (averaged over modules, and averaged over cell lines) for different sequencing feature sets (same features sets described in **Figure 2**). We observe that increase in ρ_3 penalty (which enforces the selection of the same features across all cell lines) decreases the predictive power of the model in almost all cases (no penalty (magenta curve) vs. high penalty (yellow curve)). **Supplementary Figure S9** shows the effect of ρ_1 (sparsity penalty) and ρ_2 (fused penalty) on the predictive power of the model. We observe that when using chromatin features (Motif+Chromatin, and ATAC+Motif+Chromatin), increase in ρ_1 (increasing the sparsity) and increase in ρ_2 (increasing the similarity of inferred networks) improves the predictive power of the method. Inversely, for the feature sets

that do not use chromatin (Motif, Q-Motif, and ATAC+Motif), increase in ρ_1 (sparser models) decrease the predictive power of the model. Increase in ρ_2 decrease the predictive power of the method for sparser models, while for denser networks (lower values of ρ_1) changes in ρ_2 does not significantly change the performance. Finally, we note that while ρ_1 is the main driver of sparsity of the model, increase in ρ_2 also increase the sparsity of the model (second column in **Supplementary Figure S9**).

Input datasets

We applied DRMN to study regulatory program transitions during reprogramming of differentiated cells to pluripotent cells. We obtained four datasets, one measured by array and three by sequencing. Three of these assayed mouse cell types representing different stages of reprogramming, while the fourth one focused on hepatocyte dedifferentiation. The reprogramming dataset included differentiated mouse embryonic fibroblasts (MEFs), partially reprogrammed induced pluripotent stem cells (pre-iPSCs), and pluripotent stem cells (iPSCs or embryonic stem cells, ESCs).

Reprogramming array data

We obtained measurements of gene expression and eight chromatin marks in three cell types that were processed by [19]. The measurements spanned three cell types: MEF, pre-iPSC, and iPSC. The original data were collected from multiple publications [19, 39–41]. The gene expression of 15,982 genes was measured by microarray. Eight chromatin marks were measured by ChIP-on-chip (chromatin immunoprecipitation followed by promoter microarray). For each gene promoter, each mark's value was averaged across a 8000-bp region associated with the promoter. The chromatin marks were associated with active transcription (H3K4me3, H3K9ac, H3K14ac, and H3K18ac), repression (H3K9me2, H3K9me3, H3K27me3), and transcription elongation (H3K79me2).

Reprogramming data from Chronis et al

Chronis et al. [7] assayed gene expression with RNA-seq, nine chromatin marks with ChIP-seq, and chromatin accessibility with ATAC-seq using sequencing in different stages of reprogramming (MEF, MEF48 (48 hours after start of the reprogramming process), pre-iPSC, and embryonic stem cells (ESC)). We aligned

all sequencing reads to the mouse mm9 reference genome using Bowtie2 [42]. For RNA-seq data, we quantified expression to TPMs using RSEM [43] and applied a log transform. After removing unexpressed genes ($\text{TPM} < 1$), we had 17,358 genes. The epigenomic data spanned nine histone modification marks (H3K27ac, H3K27me3, H3K36me3, H3K4me1, H3K4me2, H3K4me3, H3K79me2, H3K9ac, and H3K9me3), and chromatin accessibility (ATAC-seq). For the epigenomic datasets, we obtained per-base pair read coverage using BEDTools [44], aggregated counts within $\pm 2,500$ bps of genes' transcription start sites, and applied log transformation.

Hepatocyte dedifferentiation time course data

For the dedifferentiation time course, samples were extracted from adult mouse liver, and gene expression and chromatin accessibility was assayed by RNA-seq and ATAC-seq at 0 hours, 0.5 hours, 1 hours, 2 hours, and every 2 hours until 24 hours, and finally 36 hours (16 time points in total). All sequencing reads were aligned to mouse mm10 reference genome using Bowtie2 [42], and gene expression was quantified using RSEM [43]. Any gene with $\text{TPM} = 0$ in all time points was removed, resulting in 14,794 genes with measurement in at least one time point. Per-base pair read coverage for ATAC-seq was obtained using BEDTools [44], and counts were aggregated within $\pm 2,500$ bps of genes' transcription start sites. Both gene expression and accessibility data were quantile normalized across 16 time points and then log transformed.

Feature set generation for reprogramming datasets

We considered the following features for each gene to predict its expression. Feature datasets were processed as described above. Features involving the accessibility data were only available for the sequencing dataset (marked below with *). Any missing values in the feature data were set to 0.

- Motif network, scored by $-\log_{10}(p\text{-value})$ of motif instances (features for 353 TFs). We assembled a motif-based, context-independent network by computationally identifying links between transcription factor motifs and the genes from the expression data based on the presence of a motif instance within the gene's promoter region. We downloaded a meta-compilation of mouse motif PWMs from <http://piq.csail.mit.edu/> [45], which were sourced from multiple databases [46–48]. From the full motif list, we used only those annotated as transcription factor proteins [49]. We applied FIMO [50] to

scan the mouse genome (Ensembl version NCBIM37.66) for significant motif instances ($p < 1e - 5$). This resulted in a total of 353 TFs. Each edge in the motif network between transcription factor r and gene g is scored with the smallest p -value for any instance of r within 2kb of any TSS for g . We also assessed a z -score transformation of the p -values, which performed poorly compared to the p -value representation (not shown).

- Chromatin mark signal (feature for 8 marks in array dataset, 9 marks in sequencing dataset)
- Chromatin + motif: Concatenation of chromatin and motif features (361 features array, 362 sequencing)
- *ATAC: Log ATAC-seq signal (1 feature).
- *ATAC + motif: ATAC-seq signal and motif features (354 features)
- *Q-motif: Motif features scored by ATAC-seq signal (353 features) We quantified cell type-specific motif networks using the ATAC-seq data from the sequencing dataset. We used BedTools (`bedtools genomecov -ibam input.bam -bg -pc > output.counts`) to obtain the aggregated signal on each base pair. We defined the feature value as the log-transformed mean ATAC-seq count under each motif instance.
- *Chromatin + ATAC + motif: Concatenation of chromatin, accessibility, and motif features (363 features).
- *Chromatin + Q-motif: Concatenation of chromatin and Q-motif features (362 features)
- *Chromatin + ATAC + Q-motif: Concatenation of chromatin, accessibility, and Q-motif features (363 features).

Feature set generation for dedifferentiation datasets

We used PIQ package [45] to identify genome-wide motif instances using PWMs from CIS-BP database [51]. Motif instances were mapped to $pm2$, 500bp of genes' TSS, creating 2,856 features. Motif instances were scored by ATAC-seq signal (see Q-Motif above), and accessibility signal on promoter ($pm2$, 500bp

of genes' TSS) was added as an additional feature. Features were quantile normalized and log transformed across the 16 time points.

Experiments to evaluate DRMN and feature combinations

In our experiments, we sought to evaluate (a) whether sharing information between cell types is beneficial, and (b) which types of regulatory features are most useful to predict gene expression. We evaluated expression prediction in three-fold cross-validation, where a model was trained on two-thirds of the genes and used to predict expression for the held-aside third. We ran each experiment over a range of the number of modules, $k \in 3, 5, 7, 9, 11$.

DRMN's performance was compared to that of two baseline algorithms. *RMN* is simply DRMN run on one cell type at a time, keeping all other experimental parameters identical. *GMM* applies Gaussian Mixture Model clustering to gene expression values, and predicts the expression of a held-aside gene as the mean of the module with the highest posterior probability. We ran GMM per-cell type (*GMM-Indep*) as well as on a merged expression matrix with all three cell types (*GMM-Merged*). For GMM-Merged, the expression predictions were scored per cell-type.

We evaluated DRMN to other base line methods based on the ability to predict held-aside gene expression in three-fold cross validation. Expression prediction was evaluated using two metrics. First, the overall expression prediction was assessed using, the Pearson's correlation between actual and predicted expression of all held-aside genes. Second, the average Pearson's correlation in each module, which reflects a more fine-grained view of the utility of using regulatory features in explaining the expression levels of individual genes in a module. The first metric examines how each model (e.g., simple expression-based clustering approach versus expression and regulatory feature-based model), explains the overall variation in the data. The second metric assess the value of using additional regulatory features, such as, sequence and chromatin to predict expression. We assessed this ability using a 3-fold cross validation scheme for different values of k , number of clusters, and compared the observed (true) expression in the three test sets to predicted expression values.

Identification of transitioning genes

A transitioning gene was defined as a gene that its module assignment was changed in at least one time point/cell line. These genes were grouped into transitioning gene sets using hierarchical clustering approach (with city block as distance metric and 0.05 as distance threshold) using our in-house programs.

Prediction of regulators for transitioning genes

To identify regulators associated with transitioning gene sets in dedifferentiation dataset, we used Multi-Task Group LASSO. Briefly, in each transitioning gene sets, we solve a linear regression for each gene to predict its expression (across time points) using Q-motif features as predictor, with the group constraint to enforce selection of similar predictors for all the genes in the gene set:

$$\min_{\Theta} \|X - Y\Theta\|_2^2 + \lambda\|\Theta\|_{2,1},$$

where X is the expression matrix of genes in the transitioning gene set, and Y correspond to feature matrix of Q-motif features of regulators for genes in the transitioning gene set. We used MATLAB implementation of multi-task group LASSO from SLEPP package [52]. We performed leave-one-out cross-validation and selected regulators that were selected in at least 60% of the trained models. Additionally, we used randomized feature data to train 40 random models, and asked if the frequency of selecting a regulator was significantly higher random models (z-test with p -value < 0.05 using mean and standard deviation from the 40 random models).

Availability

The DRMN code is available at <https://github.com/Roy-lab/drmn>

Acknowledgment

This work was made possible in part by NIH NIGMS grant 1R01GM117339 to S.R. The authors thank the Center for High Throughput Computing (CHTC) for computing resources.

References

- [1] Alvaro J. Gonzalez, Manu Setty, and Christina S. Leslie. Early enhancer establishment and regulatory locus complexity shape transcriptional programs in hematopoietic differentiation. *Nature genetics*, 47(11):1249–1259, Nov 2015. 26390058[pmid].
- [2] Hatice U. Osmanbeyoglu, Fumiko Shimizu, Angela Rynne-Vidal, Petar Jelinic, Samuel C. Mok, Gabriela Chiosis, Douglas A. Levine, and Christina S. Leslie. Chromatin-informed inference of transcriptional programs in gynecologic and basal breast cancers. *bioRxiv*, 2018.
- [3] Richard A. Young. Control of the embryonic stem cell state. *Cell*, 144(6):940 – 954, 2011.
- [4] Tong Ihn Lee and Richard A Young. Transcriptional regulation and its misregulation in disease. *Cell*, 152(6):1237–1251, Mar 2013.
- [5] Joseph A Wamstad, Jeffrey M Alexander, Rebecca M Truty, Avanti Shrikumar, Fugen Li, Kirsten E Eilertson, Huiming Ding, John N Wylie, Alexander R Pico, John A Capra, Genevieve Erwin, Steven J Kattman, Gordon M Keller, Deepak Srivastava, Stuart S Levine, Katherine S Pollard, Alisha K Holloway, Laurie A Boyer, and Benoit G Bruneau. Dynamic and coordinated epigenetic regulation of developmental transitions in the cardiac lineage. *Cell*, 151:206–220, September 2012.
- [6] David Lara-Astiaso, Assaf Weiner, Erika Lorenzo-Vivas, Irina Zaretsky, Diego Adhemar Jaitin, Eyal David, Hadas Keren-Shaul, Alexander Mildner, Deborah Winter, Steffen Jung, Nir Friedman, and Ido Amit. Chromatin state dynamics during blood formation. *Science (New York, N.Y.)*, 345:943–949, August 2014.
- [7] Constantinos Chronis, Petko Fiziev, Bernadett Papp, Stefan Butz, Giancarlo Bonora, Shan Sabri, Jason Ernst, and Kathrin Plath. Cooperative binding of transcription factors orchestrates reprogramming. *Cell*, 168:442–459.e20, January 2017.
- [8] Esti Yeger-Lotem, Laura Riva, Linhui Julie Su, Aaron D Gitler, Anil G Cashikar, Oliver D King, Pavan K Auluck, Melissa L Geddie, Julie S Valastyan, David R Karger, Susan Lindquist, and Ernest

- Fraenkel. Bridging high-throughput genetic and transcriptional data reveals cellular responses to alpha-synuclein toxicity. *Nature genetics*, 41:316–323, March 2009.
- [9] Nir Yosef, Alex K. Shalek, Jellert T. Gaublomme, Hulin Jin, Youjin Lee, Amit Awasthi, Chuan Wu, Katarzyna Karwacz, Sheng Xiao, Marsela Jorgolli, David Gennert, Rahul Satija, Arvind Shakya, Diana Y. Lu, John J. Trombetta, Meenu R. Pillai, Peter J. Ratcliffe, Mathew L. Coleman, Mark Bix, Dean Tantin, Hongkun Park, Vijay K. Kuchroo, and Aviv Regev. Dynamic regulatory network controlling TH17 cell differentiation. *Nature*, 496(7446):461–468, Apr 2013.
- [10] Xianjun Dong, Melissa C Greven, Anshul Kundaje, Sarah Djebali, James B Brown, Chao Cheng, Thomas R Gingeras, Mark Gerstein, Roderic Guigó, Ewan Birney, and Zhiping Weng. Modeling gene expression using chromatin features in various cellular contexts. *Genome biology*, 13:R53, June 2012.
- [11] Thais G do Rego, Helge G Roeder, Francisco A T de Carvalho, and Ivan G Costa. Inferring epigenetic and transcriptional regulation during blood cell development with a mixture of sparse linear models. *Bioinformatics (Oxford, England)*, 28:2297–2303, September 2012.
- [12] Ankur P Parikh, Wei Wu, Ross E Curtis, and Eric P Xing. Treegl: reverse engineering tree-evolving gene networks underlying developing biological lineages. *Bioinformatics (Oxford, England)*, 27:i196–i204, July 2011.
- [13] Sushmita Roy, Margaret Werner-Washburne, and Terran Lane. A multiple network learning approach to capture system-wide condition-specific responses. *Bioinformatics*, 27(13):1832–1838, 2011.
- [14] Vladimir Jojic, Tal Shay, Katelyn Sylvia, Or Zuk, Xin Sun, Joonsoo Kang, Aviv Regev, Daphne Koller, Immunological Genome Project Consortium, Adam J. Best, Jamie Knell, Ananda Goldrath, Vladimir Jojic, Daphne Koller, Tal Shay, Aviv Regev, Nadia Cohen, Patrick Brennan, Michael Brenner, Francis Kim, Tata Nageswara Rao, Amy Wagers, Tracy Heng, Jeffrey Ericson, Katherine Rothamel, Adriana Ortiz-Lopez, Diane Mathis, Christophe Benoist, Natalie A. Bezman, Joseph C. Sun, Gundula Min-Oo, Charlie C. Kim, Lewis L. Lanier, Jennifer Miller, Brian Brown, Miriam Merad, Emmanuel L. Gautier, Claudia Jakubzick, Gwendalyn J. Randolph, Paul Monach, David A. Blair, Michael L. Dustin, Susan A. Shinton, Richard R. Hardy, David Laidlaw, Jim Collins, Roi Gazit, Derrick J. Rossi, Nidhi

- Malhotra, Katelyn Sylvia, Joonsoo Kang, Taras Kreslavsky, Anne Fletcher, Kutlu Elpek, Angelique Bellemarte-Pelletier, Deepali Malhotra, and Shannon Turley. Identification of transcriptional regulators in the mouse immune system. *Nat Immunol*, 14(6):633–643, Jun 2013.
- [15] Wuming Gong, Naoko Koyano-Nakagawa, Tongbin Li, and Daniel J Garry. Inferring dynamic gene regulatory networks in cardiac differentiation through the integration of multi-dimensional data. *BMC bioinformatics*, 16:74, March 2015.
- [16] Emma Pierson, GTEx Consortium, Daphne Koller, Alexis Battle, Sara Mostafavi, Kristin G Ardlie, Gad Getz, Fred A Wright, Manolis Kellis, Simona Volpi, and Emmanouil T Dermizakis. Sharing and specificity of co-expression networks across 35 human tissues. *PLoS Comput Biol*, 11(5):e1004220, May 2015.
- [17] Christopher Koch, Jay Konieczka, Toni Delorey, Ana Lyons, Amanda Socha, Kathleen Davis, Sara A. Knaack, Dawn Thompson, Erin K. O’Shea, Aviv Regev, and Sushmita Roy. Inference and evolutionary analysis of genome-scale regulatory networks in large phylogenies. *Cell Systems*, 4(5):543 – 558.e8, 2017.
- [18] Jason Ernst, Oded Vainas, Christopher T Harbison, Itamar Simon, and Ziv Bar-Joseph. Reconstructing dynamic regulatory maps. *Molecular systems biology*, 3:74, 2007.
- [19] Sushmita Roy and Rupa Sridharan. Chromatin module inference on cellular trajectories identifies key transition points and poised epigenetic states in diverse developmental processes. *Genome research*, 27:1250–1262, July 2017.
- [20] Eran Segal, Michael Shapira, Aviv Regev, Dana Pe’er, David Botstein, Daphne Koller, and Nir Friedman. Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data. *Nat. Genet.*, 34(2):166–176, May 2003.
- [21] Su-In Lee, Aimée M. Dudley, David Drubin, Pamela A. Silver, Nevan J. Krogan, Dana Pe’er, and Daphne Koller. Learning a prior on regulatory potential from eQTL data. *PLoS Genet*, 5(1):e1000358+, January 2009.

- [22] R. Gherzi, M. Trabucchi, M. Ponassi, I.-E. Gallouzi, M. G. Rosenfeld, and P. Briata. Akt2-mediated phosphorylation of pitx2 controls ccnd1 mrna decay during muscle cell differentiation. *Cell Death & Differentiation*, 17(6):975–983, Jun 2010.
- [23] Anna B. Osipovich, Qiaoming Long, Elisabetta Manduchi, Rama Gangula, Susan B. Hipkens, Judsen Schneider, Tadashi Okubo, Christian J. Stoeckert, Shinji Takada, and Mark A. Magnuson. Insm1 promotes endocrine cell differentiation by modulating the expression of a network of genes that includes neurog3 and ripply3. *Development*, 141(15):2939–2949, 2014.
- [24] MF Cole, SE Johnstone, JJ Newman, MH Kagey, and RAE Young. Tcf3 is an integral component of the core regulatory circuitry of embryonic stem cells. *Genes & development*, 22(6):746–755, Mar 2008.
- [25] Frederic Lluis, Luigi Ombrato, Elisa Pedone, Stefano Pepe, Bradley J. Merrill, and Maria Pia Cosma. T-cell factor 3 (tcf3) deletion increases somatic cell reprogramming by inducing epigenome modifications. *Proceedings of the National Academy of Sciences*, 108(29):11912–11917, 2011.
- [26] Greetje Elaut, Tom Henkens, Peggy Papeleu, Sarah Snykers, Mathieu Vinken, Tamara Vanhaecke, and Vera Rogiers. Molecular mechanisms underlying the dedifferentiation process of isolated hepatocytes and their cultures. *Current Drug Metabolism*, 7(6):629–660, 2006.
- [27] Manman Wang, Ling Tang, Dahai Liu, Qi-Long Ying, and Shoudong Ye. The transcription factor gbx2 induces expression of kruppel-like factor 4 to maintain and induce nave pluripotency of embryonic stem cells. *Journal of Biological Chemistry*, 292(41):17121–17128, 2017.
- [28] Iliaria Laudadio, Isabelle Manfroid, Younes Achouri, Dominic Schmidt, Michael D. Wilson, Sabine Cordi, Lieven Thorrez, Laurent Knoops, Patrick Jacquemin, Frans Schuit, Christophe E. Pierreux, Duncan T. Odom, Bernard Peers, and Frederic P. Lemaigre. A feedback loop between the liver-enriched transcription factor network and mir-122 controls hepatocyte differentiation. *Gastroenterology*, 142(1):119–129, Jan 2012.
- [29] Avinash Thakur, Jasper C.H. Wong, Evan Y. Wang, Jeremy Lotto, Donghwan Kim, Jung-Chien Cheng, Matthew Mingay, Rebecca Cullum, Vaishali Moudgil, Nafeel Ahmed, Shu-Huei Tsai, Wei

- Wei, Colum P. Walsh, Tabea Stephan, Misha Bilenky, Bettina M. Fuglerud, Mohammad M. Karimi, Frank J. Gonzalez, Martin Hirst, and Pamela A. Hoodless. Hepatocyte nuclear factor 4-alpha is essential for the active epigenetic state at enhancers in mouse liver. *Hepatology*, 70(4):1360–1376, 2019.
- [30] Wenjie Yu, Xiao Li, Steven Eliason, Miguel Romero-Bustillos, Ryan J. Ries, Huojun Cao, and Brad A. Amendt. Irx1 regulates dental outer enamel epithelial and lung alveolar type ii epithelial differentiation. *Developmental biology*, 429(1):44–55, Sep 2017. 28746823[pmid].
- [31] Chih-I Tai and Qi-Long Ying. Gbx2, a lif/stat3 target, promotes reprogramming to and retention of the pluripotent ground state. *Journal of Cell Science*, 126(5):1093–1098, 2013.
- [32] Raven Diacou, Yilin Zhao, Deyou Zheng, Ales Cvekl, and Wei Liu. Six3 and six6 are jointly required for the maintenance of multipotent retinal progenitors through both positive and negative regulation. *Cell Reports*, 25(9):2510 – 2523.e4, 2018.
- [33] Silvia Buonamici, Soumen Chakraborty, Vitalyi Senyuk, and Giuseppina Nucifora. The role of evi1 in normal and leukemic cells. *Blood Cells, Molecules, and Diseases*, 31(2):206 – 212, 2003.
- [34] Jian Ming Khor, Jennifer Guerrero-Santoro, and Charles A. Etensohn. Genome-wide identification of binding sites and gene targets of alx1, a pivotal regulator of echinoderm skeletogenesis. *Development*, 146(16), 2019.
- [35] Manu Setty and Christina S. Leslie. Seqgl identifies context-dependent binding signals in genome-wide regulatory element maps. *PLOS Computational Biology*, 11(5):1–21, 05 2015.
- [36] Jiayu Zhou, Jianhui Chen, and Jieping Ye. Malsar: Multi-task learning via structural regularization – users manual version 1.1, 2012.
- [37] Yu. Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1):127–152, May 2005.
- [38] Yu. Nesterov. Gradient methods for minimizing composite objective function. CORE Discussion Papers 2007076, Universit catholique de Louvain, Center for Operations Research and Econometrics (CORE), 2007.

- [39] Nimet Maherali, Rupa Sridharan, Wei Xie, Jochen Utikal, Sarah Eminli, Katrin Arnold, Matthias Stadtfeld, Robin Yachechko, Jason Tchieu, Rudolf Jaenisch, Kathrin Plath, and Konrad Hochedlinger. Directly reprogrammed fibroblasts show global epigenetic remodeling and widespread tissue contribution. *Cell stem cell*, 1:55–70, June 2007.
- [40] Rupa Sridharan, Jason Tchieu, Mike J Mason, Robin Yachechko, Edward Kuoy, Steve Horvath, Qing Zhou, and Kathrin Plath. Role of the murine reprogramming factors in the induction of pluripotency. *Cell*, 136:364–377, January 2009.
- [41] Rupa Sridharan, Michelle Gonzales-Cope, Constantinos Chronis, Giancarlo Bonora, Robin McKee, Chengyang Huang, Sanjeet Patel, David Lopez, Nilamadhab Mishra, Matteo Pellegrini, Michael Carey, Benjamin A Garcia, and Kathrin Plath. Proteomic and genomic approaches reveal critical functions of H3K9 methylation and heterochromatin protein-1 γ in reprogramming to pluripotency. *Nature cell biology*, 15:872–882, July 2013.
- [42] Ben Langmead and Steven L Salzberg. Fast gapped-read alignment with bowtie 2. *Nature methods*, 9:357–359, March 2012.
- [43] Bo Li and Colin N Dewey. Rsem: accurate transcript quantification from rna-seq data with or without a reference genome. *BMC bioinformatics*, 12:323, August 2011.
- [44] Aaron R. Quinlan and Ira M. Hall. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, 26(6):841–842, 01 2010.
- [45] Richard I. Sherwood, Tatsunori Hashimoto, Charles W. O’Donnell, Sophia Lewis, Amira A. Barkal, John Peter van Hoff, Vivek Karun, Tommi Jaakkola, and David K. Gifford. Discovery of directional and nondirectional pioneer transcription factors by modeling dnase profile magnitude and shape. *Nat Biotechnol*, 32(2):171–178, Feb 2014.
- [46] V Matys, E Fricke, R Geffers, E Gössling, M Haubrock, R Hehl, K Hornischer, D Karas, A E Kel, O V Kel-Margoulis, D-U Kloos, S Land, B Lewicki-Potapov, H Michael, R Münch, I Reuter, S Rotert, H Saxel, M Scheer, S Thiele, and E Wingender. Transfac: transcriptional regulation, from patterns to profiles. *Nucleic acids research*, 31:374–378, January 2003.

- [47] Albin Sandelin, Wynand Alkema, Pär Engström, Wyeth W Wasserman, and Boris Lenhard. Jasparr: an open-access database for eukaryotic transcription factor binding profiles. *Nucleic acids research*, 32:D91–D94, January 2004.
- [48] Michael F Berger, Anthony A Philippakis, Aaron M Qureshi, Fangxue S He, Preston W Estep, and Martha L Bulyk. Compact, universal dna microarrays to comprehensively determine transcription-factor binding site specificities. *Nature biotechnology*, 24:1429–1435, November 2006.
- [49] Timothy Ravasi, Harukazu Suzuki, Carlo V. Cannistraci, Shintaro Katayama, Vladimir B. Bajic, Kai Tan, Altuna Akalin, Sebastian Schmeier, Mutsumi Kanamori-Katayama, Nicolas Bertin, Piero Carninci, Carsten O. Daub, Alistair R. R. Forrest, Julian Gough, Sean Grimmond, Jung-Hoon Han, Takehiro Hashimoto, Winston Hide, Oliver Hofmann, Atanas Kamburov, Mandeep Kaur, Hideya Kawaji, Atsuta Kubosaki, Timo Lassmann, Erik van Nimwegen, Cameron R. MacPherson, Chihiro Ogawa, Aleksandar Radovanovic, Ariel Schwartz, Rohan D. Teasdale, Jesper Tegnér, Boris Lenhard, Sarah A. Teichmann, Takahiro Arakawa, Noriko Ninomiya, Kayoko Murakami, Michihira Tagami, Shiro Fukuda, Kengo Imamura, Chikatoshi Kai, Ryoko Ishihara, Yayoi Kitazume, Jun Kawai, David A. Hume, Trey Ideker, and Yoshihide Hayashizaki. An atlas of combinatorial transcriptional regulation in mouse and man. *Cell*, 140(5):744–752, March 2010.
- [50] Charles E. Grant, Timothy L. Bailey, and William Stafford Noble. Fimo: scanning for occurrences of a given motif. *Bioinformatics*, 27(7):1017–1018, Apr 2011.
- [51] Matthew T. Weirauch, Ally Yang, Mihai Albu, Atina G. Cote, Alejandro Montenegro-Montero, Philipp Drewe, Hamed S. Najafabadi, Samuel A. Lambert, Ishminder Mann, Kate Cook, Hong Zheng, Alejandra Goity, Harm van Bakel, Jean-Claude Lozano, Mary Galli, Mathew G. Lewsey, Eryong Huang, Tuhin Mukherjee, Xiaoting Chen, John S. Reece-Hoyes, Sridhar Govindarajan, Gad Shaulsky, Albertha J M. Walhout, François-Yves Bouget, Gunnar Ratsch, Luis F. Larrondo, Joseph R. Ecker, and Timothy R. Hughes. Determination and inference of eukaryotic transcription factor sequence specificity. *Cell*, 158(6):1431–1443, Sep 2014.
- [52] Jolanta Jura, Paulina Wegrzyn, Michal Korostynski, Krzysztof Guzik, Malgorzata Oczko-Wojciechowska, Michal Jarzab, Malgorzata Kowalska, Marcin Piechota, Ryszard Przewlocki, and

Aleksander Koj. Identification of interleukin-1 and interleukin-6-responsive genes in human monocyte-derived macrophages using microarrays. *Biochimica et Biophysica Acta (BBA) - Gene Regulatory Mechanisms*, 1779(6):383 – 389, 2008.

Main Figures

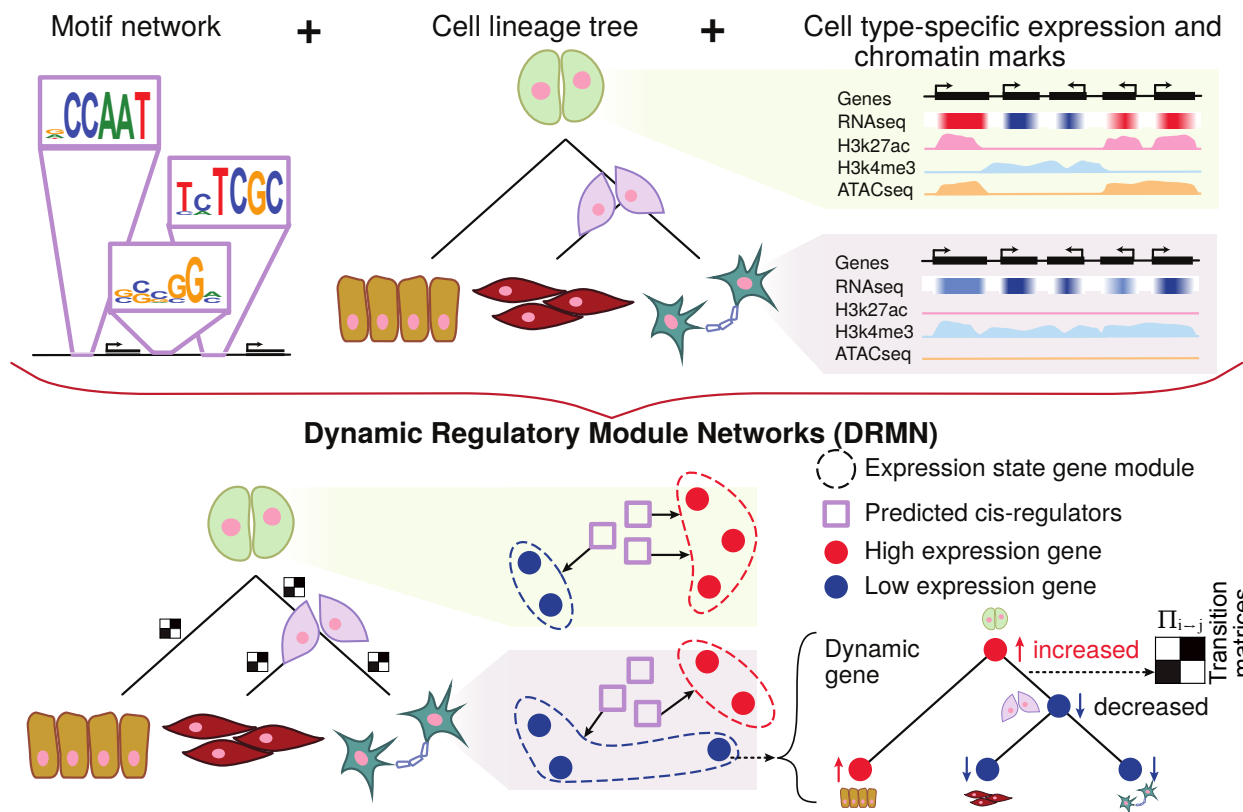


Figure 1: The outline of Dynamic Regulatory Module Networks (DRMN) method. Inputs are a lineage tree over the cell types, cell type-specific expression levels, a shared skeleton regulatory network (*e.g.* sequence specific motif network), and optionally cell type-specific features such as histone modification marks or chromatin accessibility signal. The output is a learned DRMN, which consists of cell type-specific expression state modules, their regulatory programs, and transition matrices describing dynamics between the cell types.

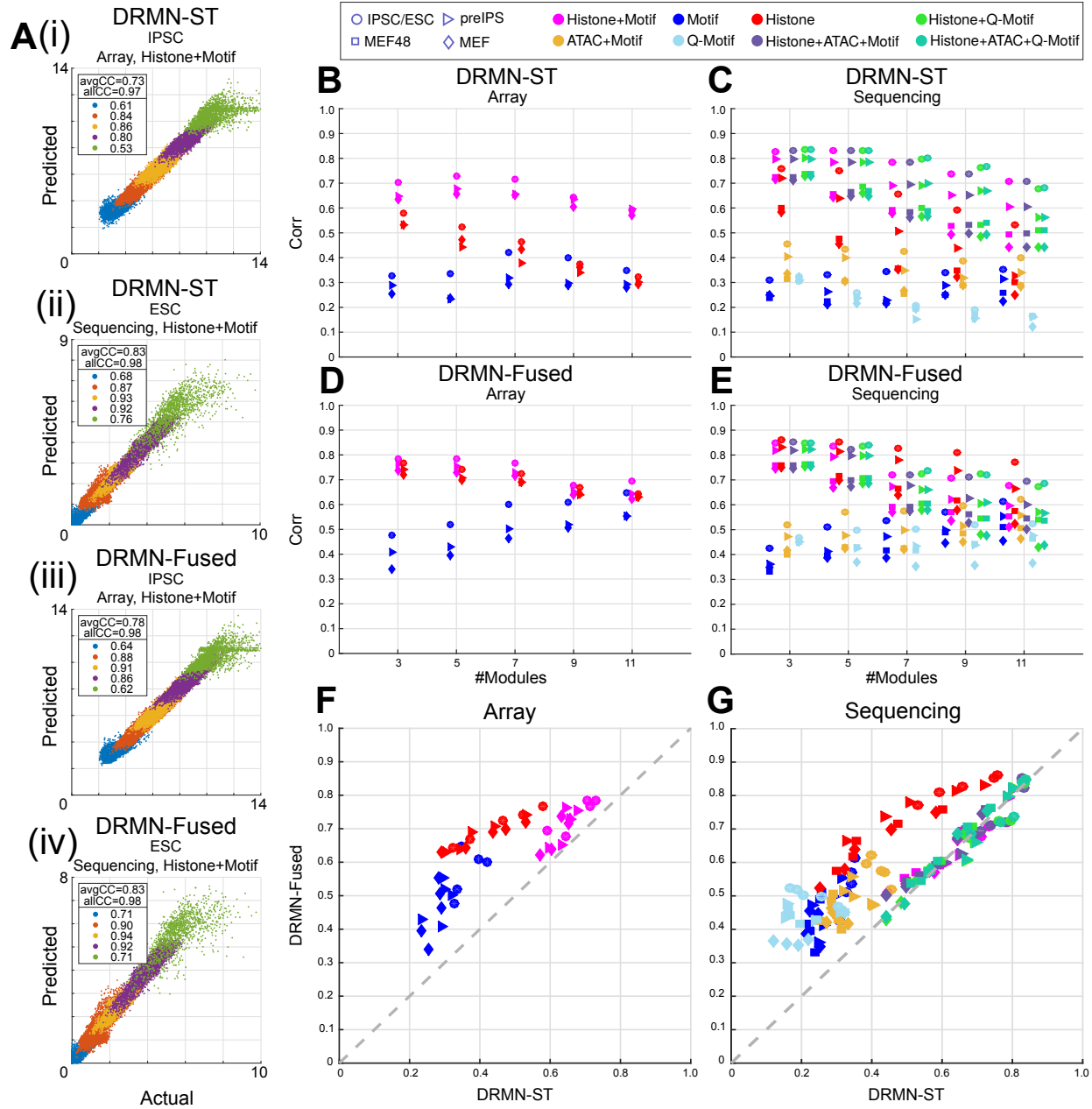


Figure 2: **(A)** Predicted expression *vs.* observed expression, for iPSC/ESC, for **(i)** DRMN-ST on array dataset, **(ii)** DRMN-ST on sequencing dataset, **(iii)** DRMN-FUSED on array dataset, and **(iv)** DRMN-FUSED on sequencing dataset. Average per-module correlation for individual cell lines as a function of different number of modules, for **(B)** DRMN-ST on array dataset, **(C)** DRMN-ST on sequencing dataset, **(D)** DRMN-FUSED on array dataset, and **(E)** DRMN-FUSED on sequencing dataset. Average per-module correlation for individual cell lines for DRMN-ST *vs.* DRMN-FUSED for **(F)** array dataset, and **(G)** sequencing dataset. Each shape correspond to a cell line and each color correspond to a different feature set.

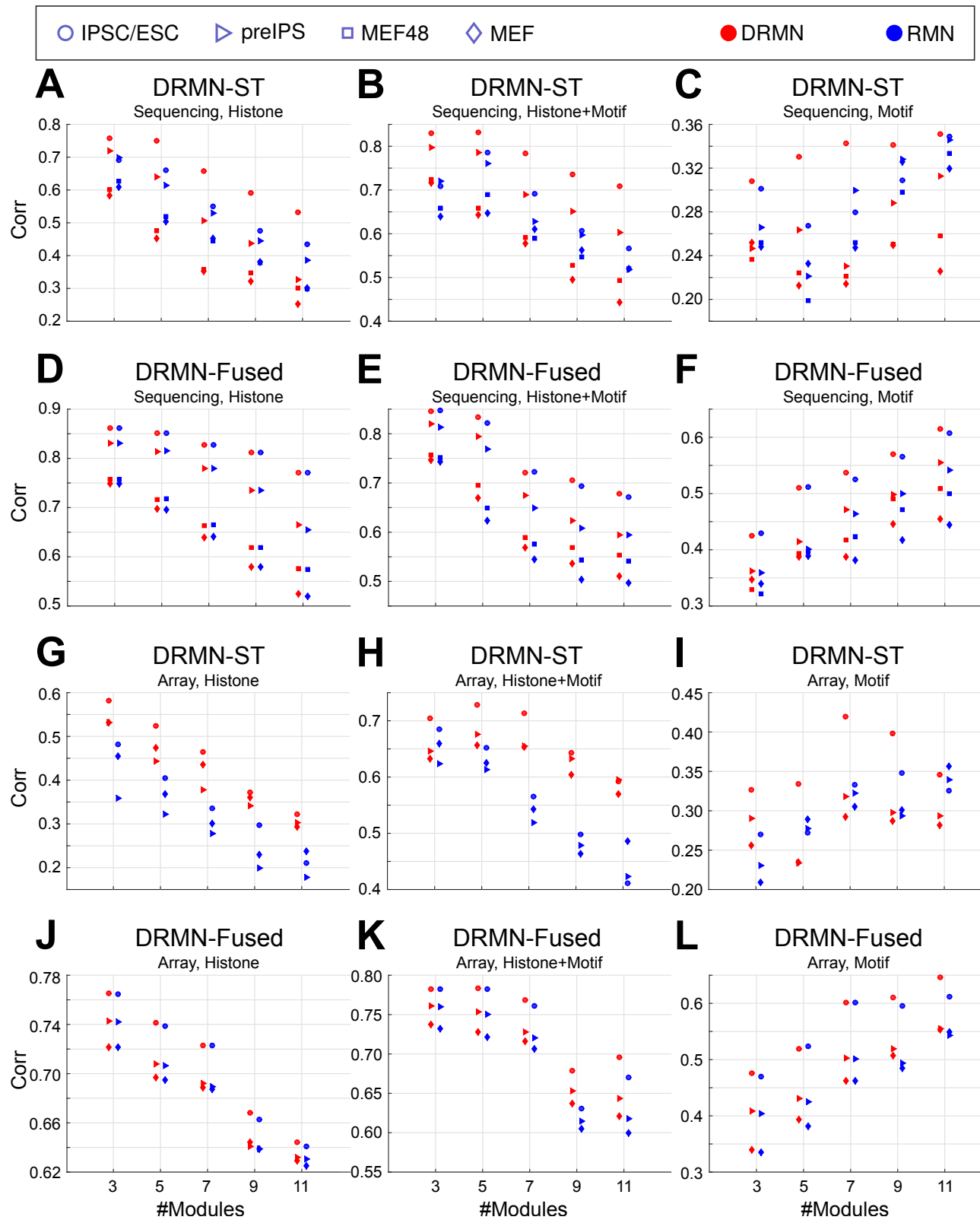


Figure 3: Average per-module correlation for individual cell lines as a function of different number of modules for single task and multi task versions of the method, for **A-C)** DRMN-ST on sequencing dataset, **D-F)** DRMN-FUSED on sequencing dataset, **G-I)** DRMN-ST on array dataset, and **J-L)** DRMN-FUSED on array dataset. Each shape correspond to a cell line and each color correspond to a different method.

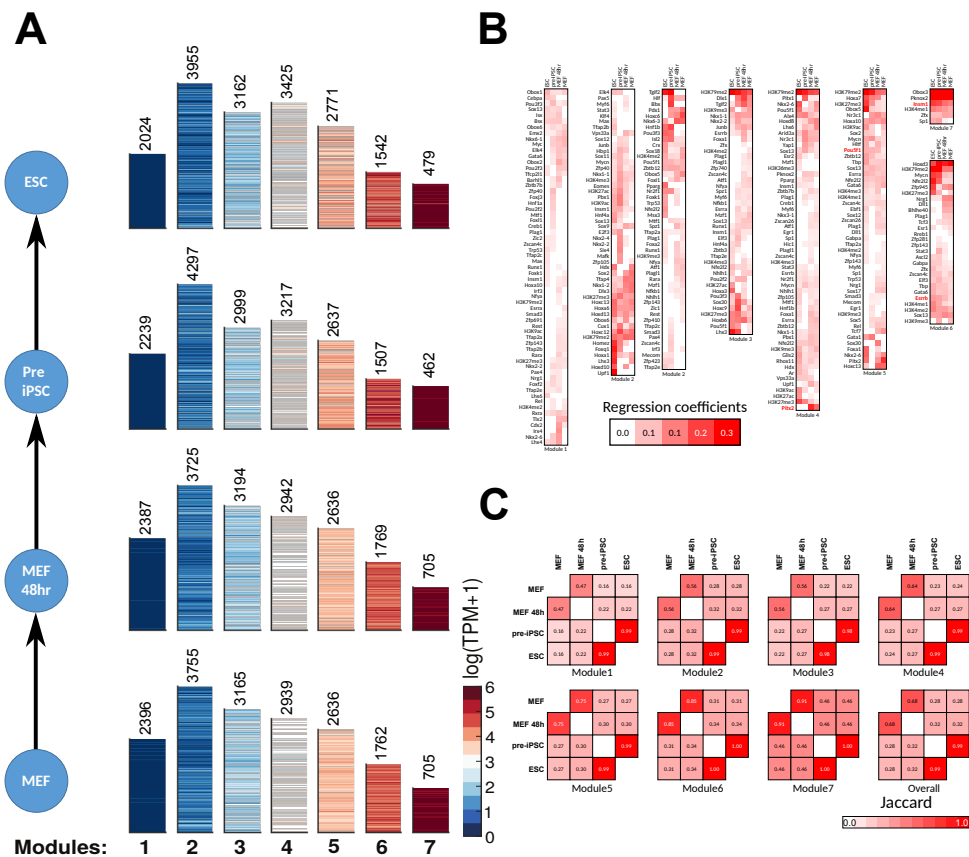


Figure 4: Annotation of DRMN modules, for models trained on chromatin+q-motif with $k = 7$ modules, for sequencing dataset. **(A)** The gene expression pattern of the 7 modules. The number above the heatmap correspond to the number of genes in that module. **(B)** Inferred regulatory program for different modules and cell lines. **(C)** Similarity of modules across cell lines.

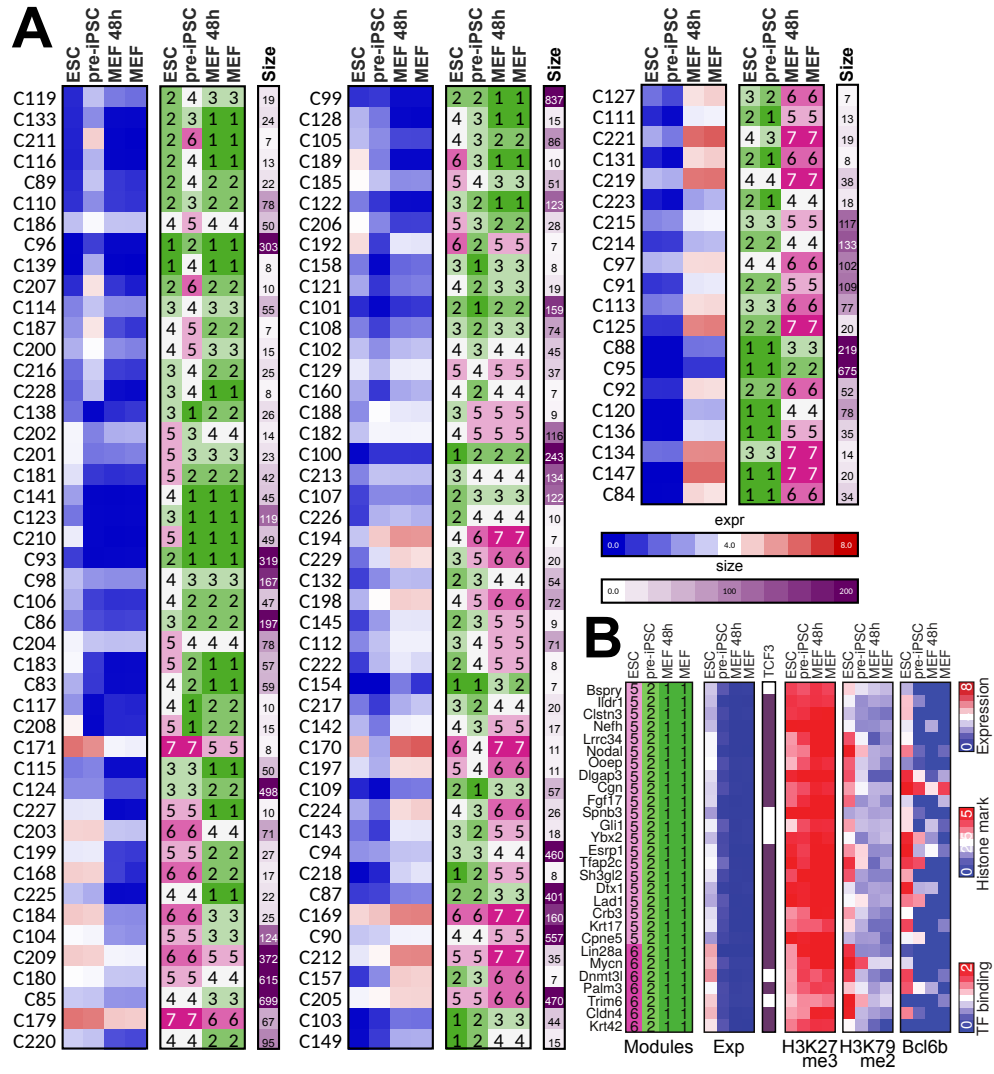


Figure 5: **(A)** Bird's eye view of groups of transitioning genes. On left is the average expression of the genes in the group, in the middle is module assignments, and on right the number of genes in that group. **(B)** Module assignments, gene expression, presence or absence of TCF3 motif, and feature values for H3K27me3, H3K79me2, and Bcl6b, for genes that transition from low to high expression between differentiated or partially reprogrammed cells and ESCs/iPSC.

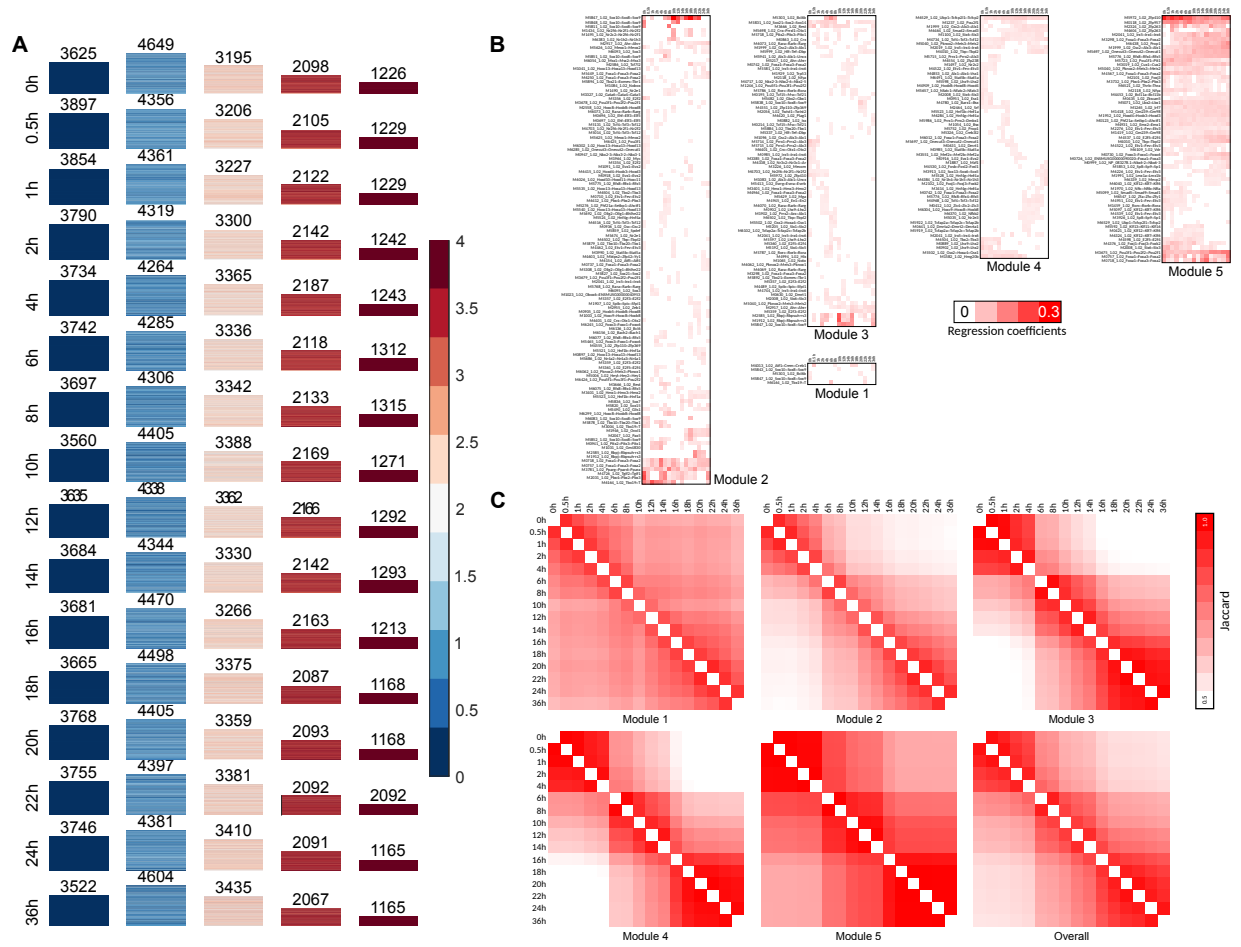


Figure 6: Annotation of DRMN modules for dedifferentiation dataset with $k = 5$ modules. **(A)** The gene expression pattern of the 5 modules. The number above the heatmap correspond to the number of genes in that module. **(B)** Inferred regulatory program for different modules and time points. **(C)** Similarity of modules across time point.

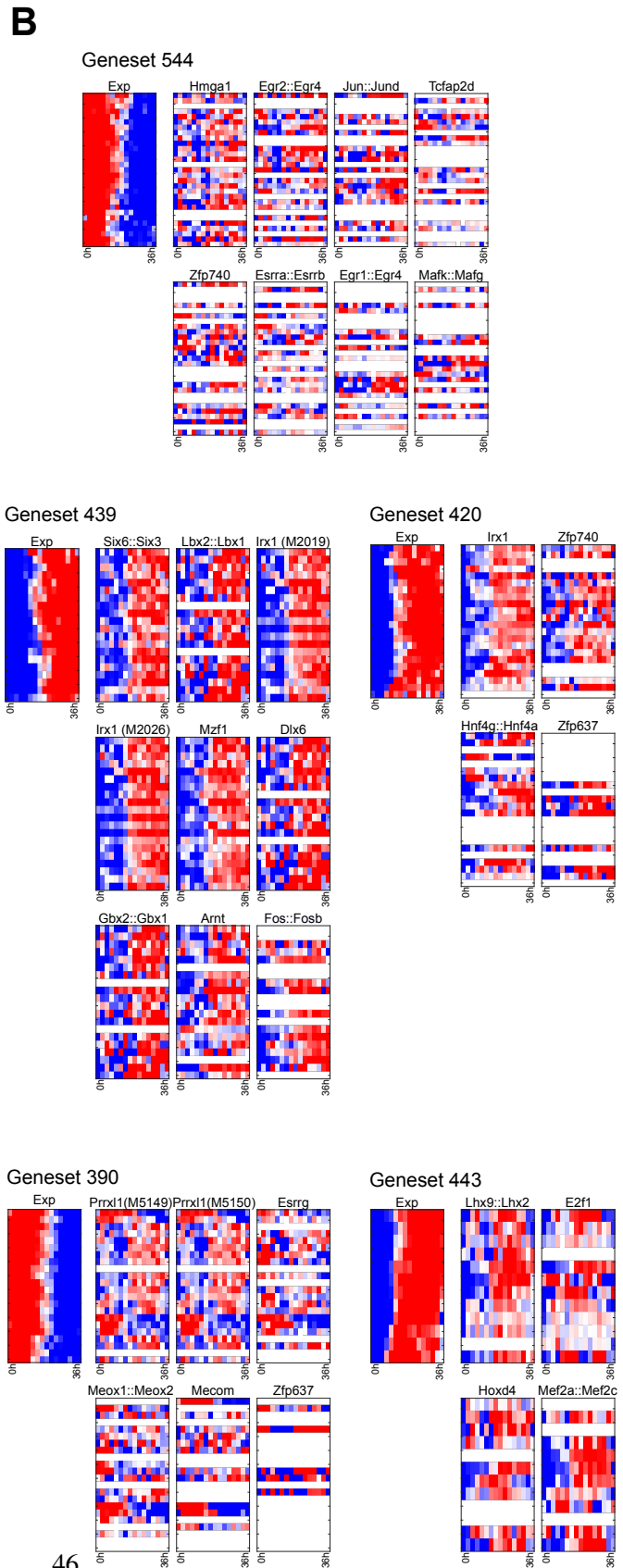
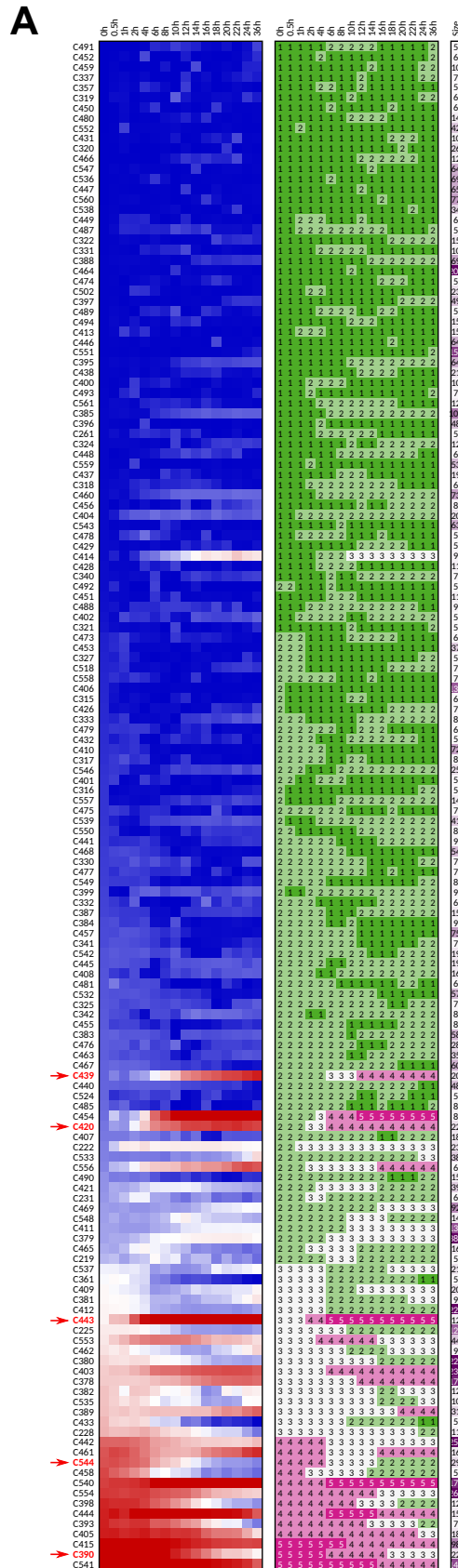


Figure 7: **(A)** Bird's eye view of groups of transitioning genes. On left is the average expression of the genes in the group, in the middle is module assignments, and on right the number of genes in that group. **(B)** Regulators associated some of transitioning gene sets.

Supplementary Figures

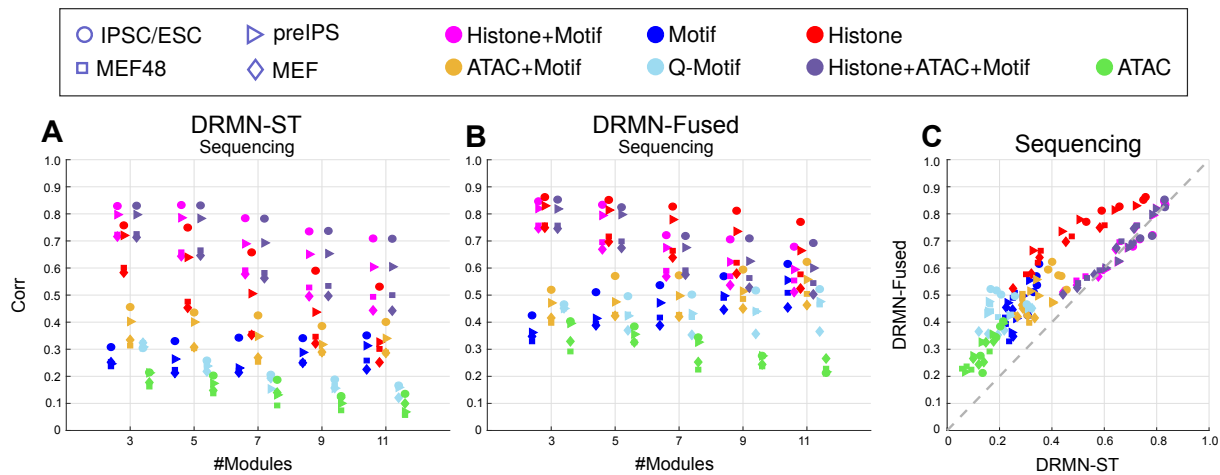


Figure S1: Performance of ATAC as a single feature (green) *vs.* Motif alone (blue), Histone+Motif (magenta), Q-Motif (light blue), ATAC+Motif (yellow), Histone alone (red), and Histone+ATAC+Motif (dark purple) using **A**) DRMN-ST and **B**) DRMN-Fused. **C**) Shows the performance of DRMN-ST *vs.* DRMN-Fused for the same features.

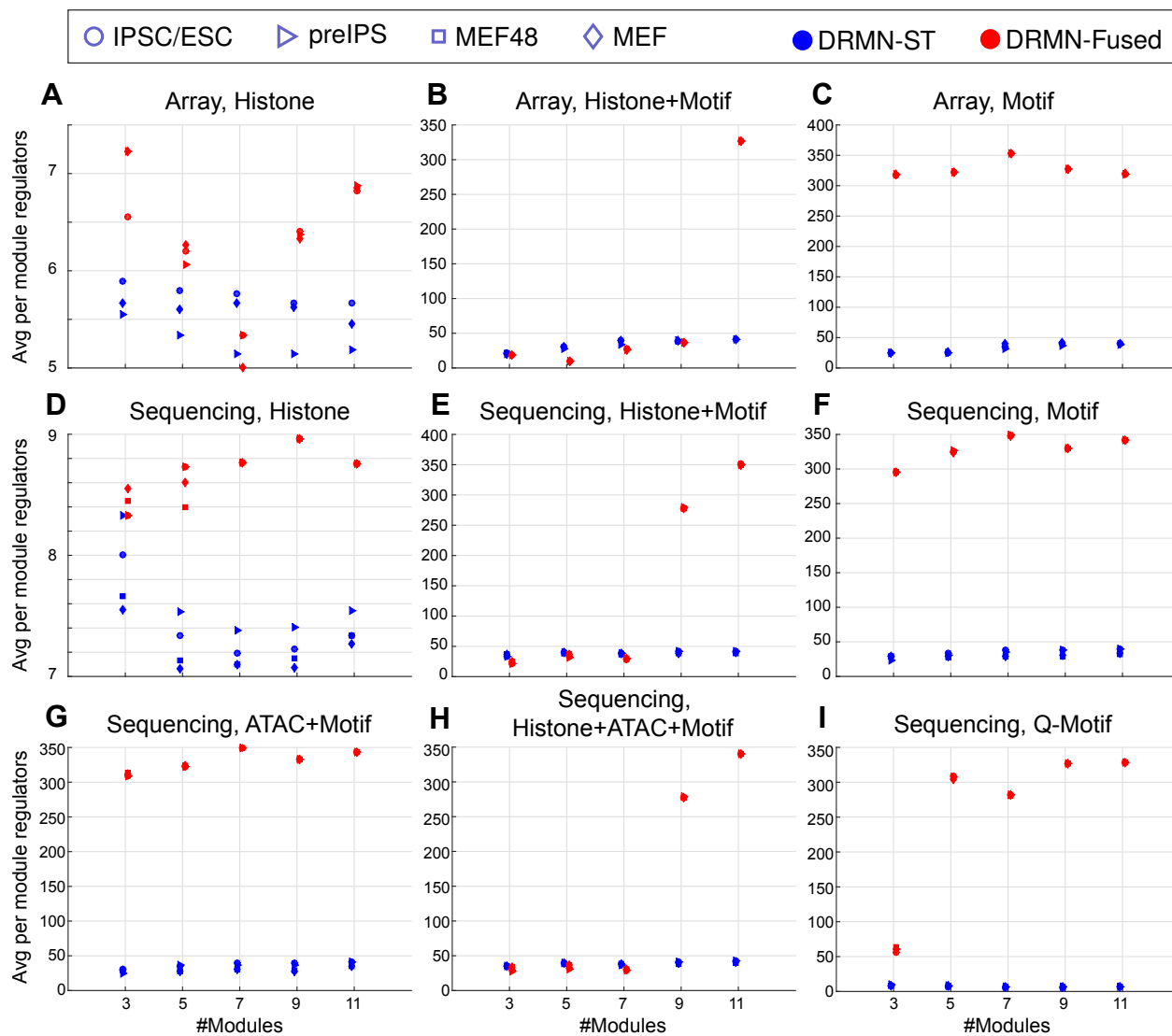


Figure S2: Comparing average number of regulators per module between DRMN ST and Fused models for different dataset/feature combinations **A-C**) feature sets in array dataset, and **D-I**) feature sets in sequencing dataset.

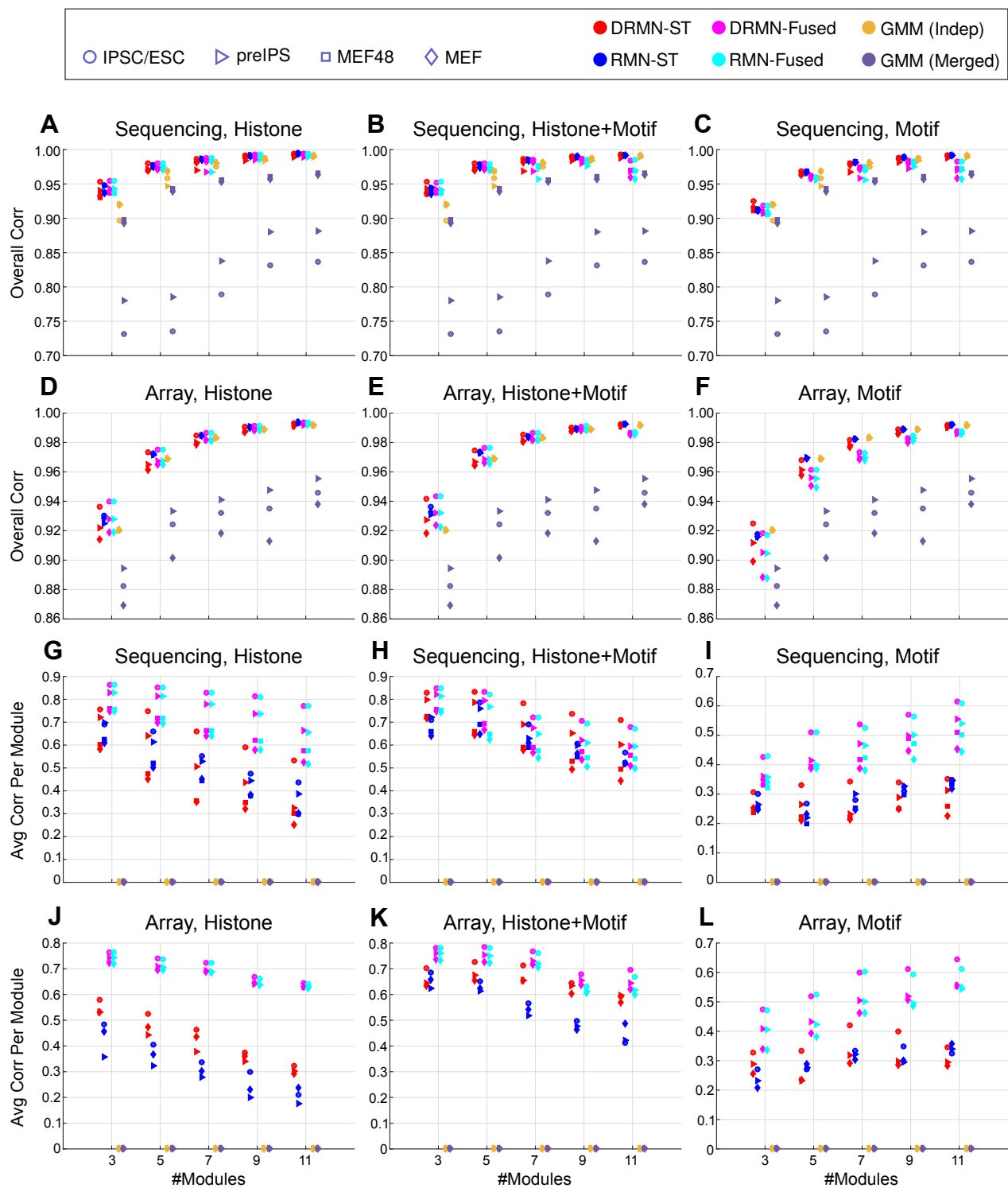


Figure S3: Comparing DRMN's expression predictions to baseline approaches. Each panel compares six algorithms (legend) on the basis of one correlation metric (y-axis) across a range of k (x-axis). Each dot represents the results for one cell type averaged over 3 fold of cross validation. The columns show results for chromatin mark features (left), combination of chromatin and motif (middle), and the motif features (right). **(A-F)** Comparison based on Pearson correlation of predicted to true expression for held-aside genes. **(G-L)** Comparison of per-module correlation coefficients, averaged across k modules.

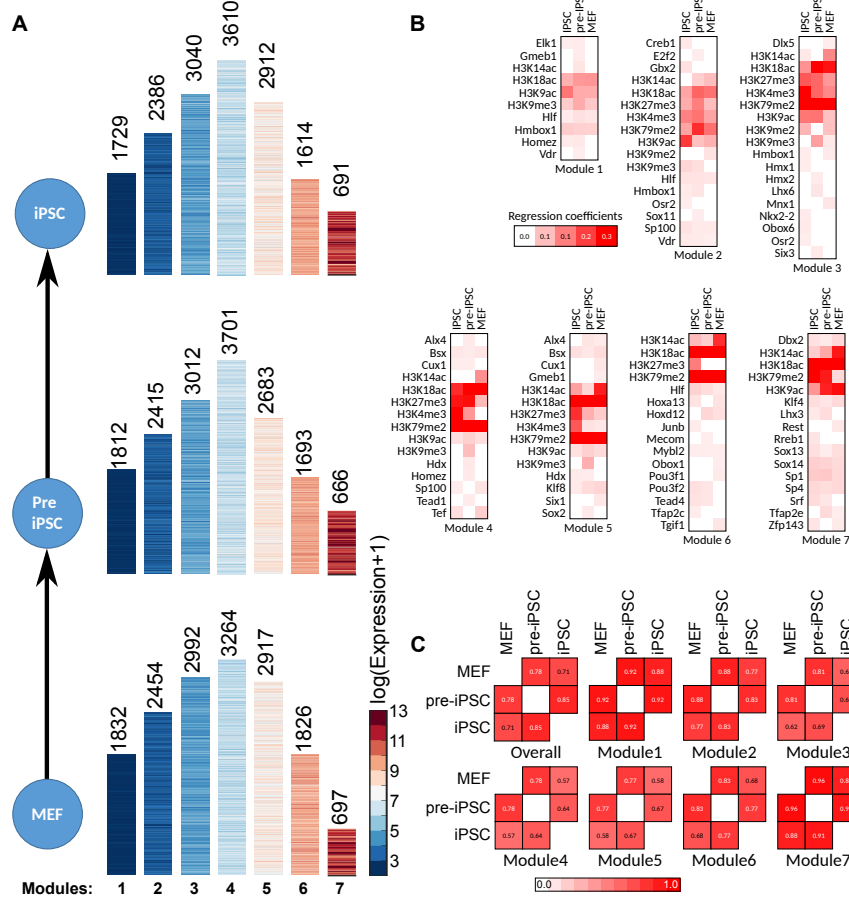


Figure S4: Annotation of DRMNs modules, for models trained on chromatin+motif with $k = 7$ modules, for array dataset. **(A)** The gene expression pattern of the 7 modules. The number above the heatmap correspond to the number of genes in that module. **(B)** Inferred regulatory program for different modules and cell lines. **(C)** Similarity of modules across cell lines.

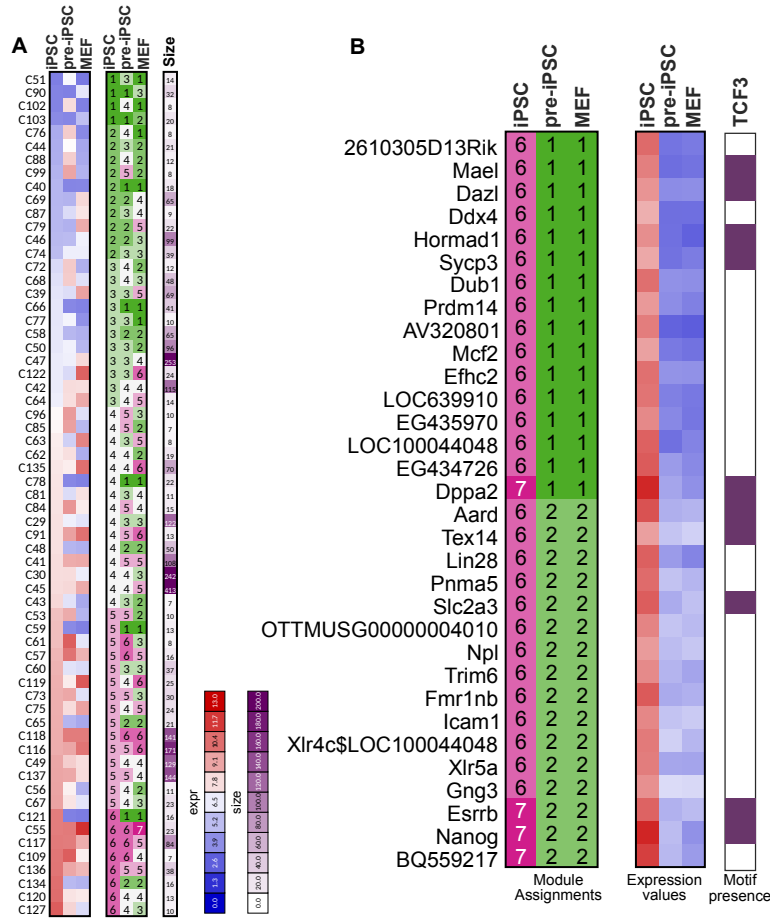


Figure S5: **(A)** Bird's eye view of groups of transitioning genes; on left is the average expression of the genes in the group, in the middle is module assignments, and on right the number of genes in that group. **(B)** Module assignments (left) and gene expression (middle), and presence or absence of TCF3 motif for genes that transition from low to high expression between differentiated or partially reprogrammed cells and ESCs/iPSC.



Figure S6: Enriched GO terms in dedifferentiation dataset that show significant changed in enrichment between earlier time points and later time points.

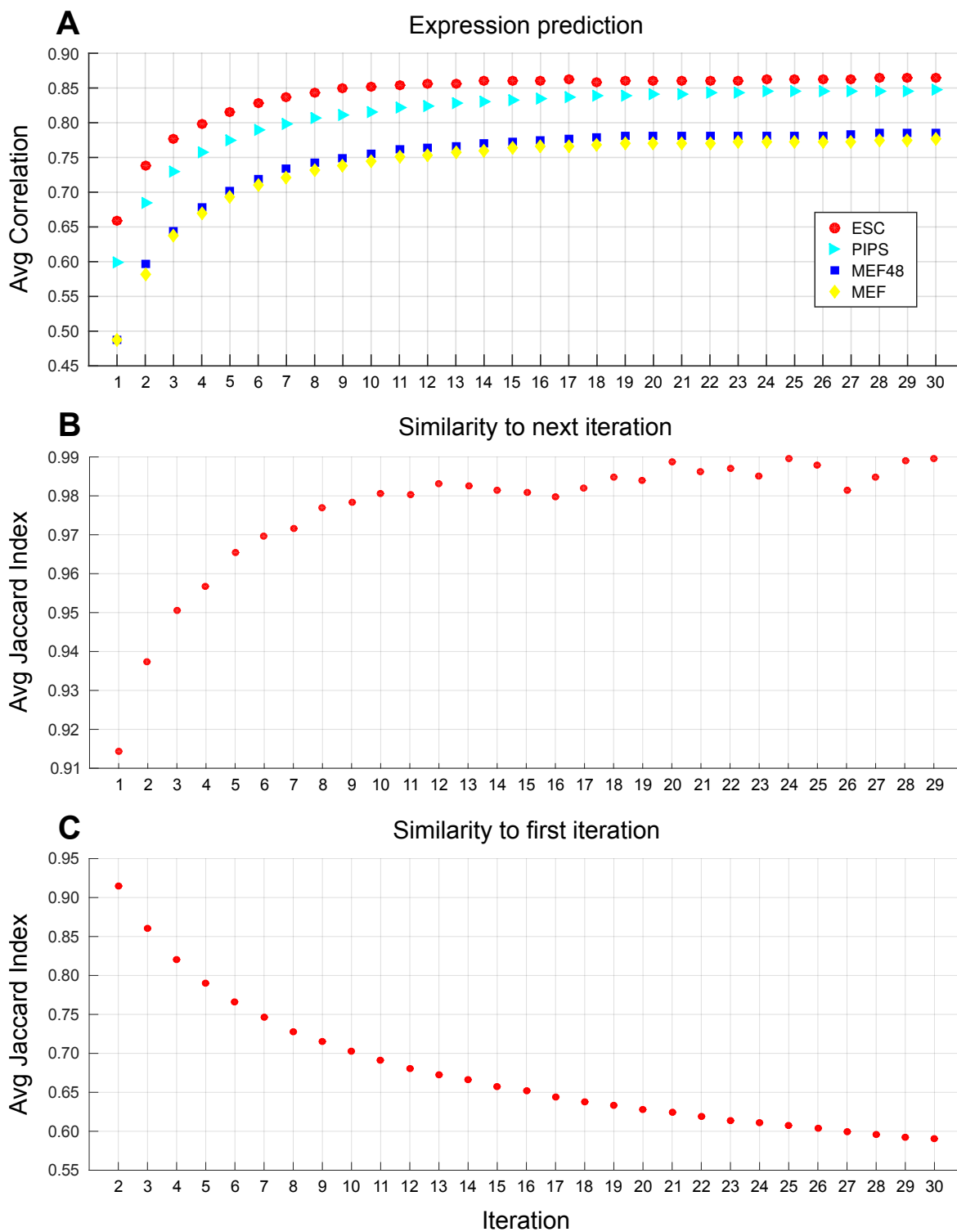


Figure S7: The effect of iteration on performance of DRMN (on Motif+Chromatin, $k=3$, using fused LASSO ($\rho_1 = 100, \rho_2 = 50, \rho_3 = 0$)). **A)** Average correlation (over all modules) for each cell line, as a function of number of iterations. Each marker corresponds to a cell line. **B)** Similarity of module assignments between consecutive iterations. **C)** Similarity of module assignments between iteration 1 and iteration i .

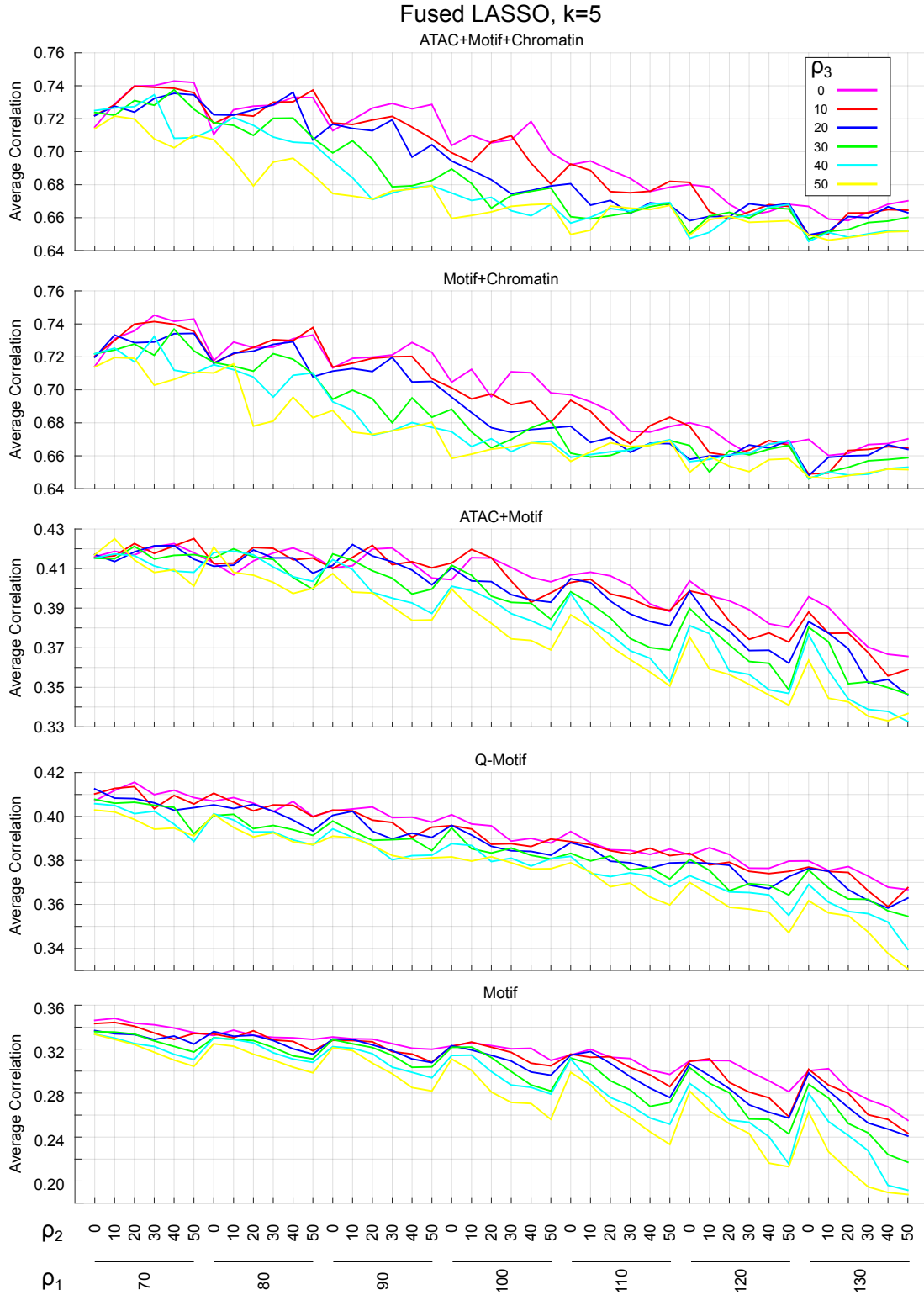


Figure S8: The effect of group penalty on DRMN performance. Each panel corresponds to one of sequencing feature sets. In each panel, different colors corresponds to different values of ρ_3 , the group LASSO penalty (enforcing selection of same features for all cell lines).

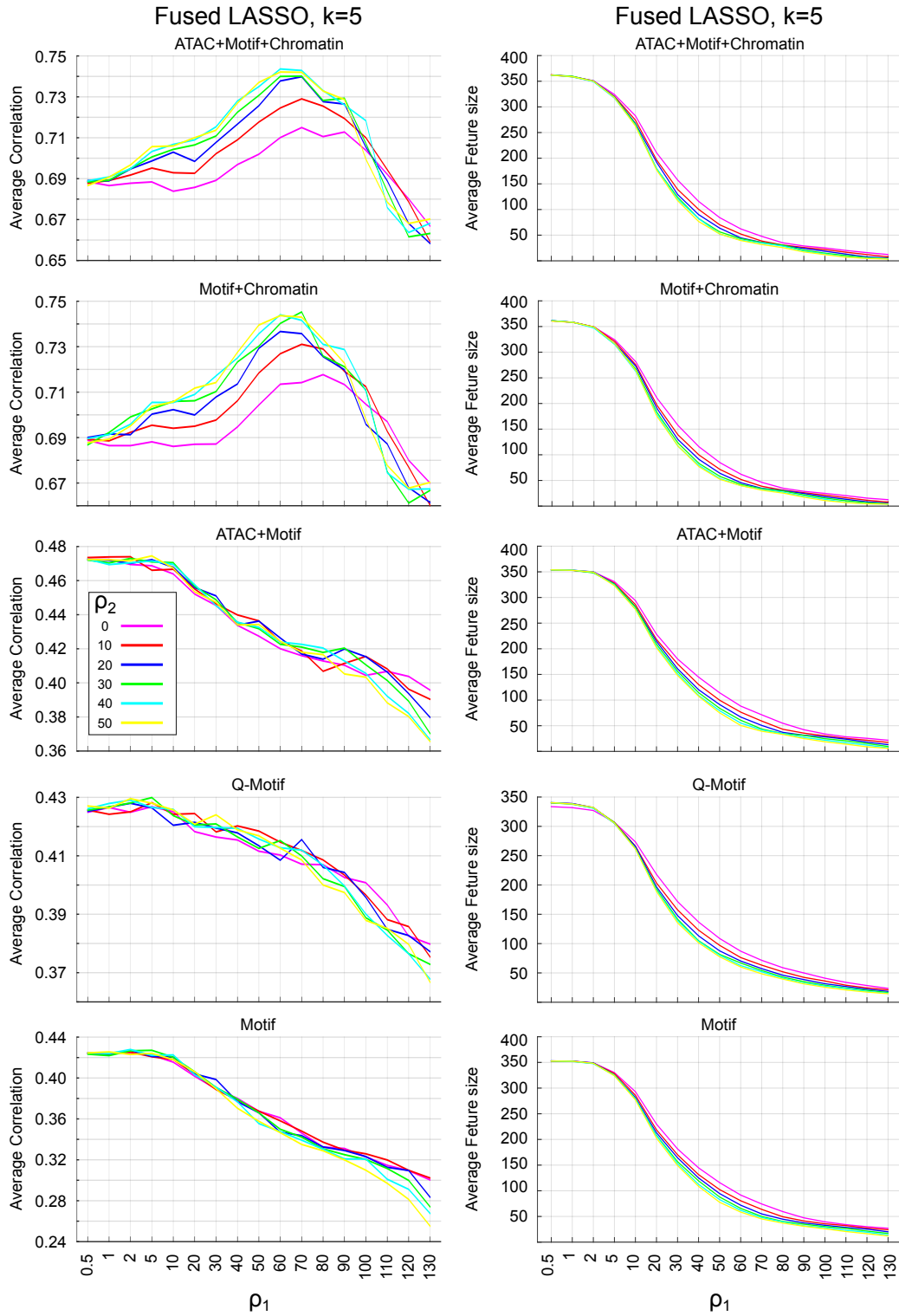


Figure S9: The effect of fused penalty on DRMN performance and the number of selected features. Each panel corresponds to one of sequencing feature sets. In each panel, different colors correspond to different values of ρ_2 , the fused LASSO penalty (enforcing similarity of features selected for closer cell lines). First column shows the average correlation of predicted expression, and second columns shows the average number of selected features (averaged over cell lines and modules).