

1 Introduction

One of the central cellular processes underlying development is transcriptional regulation. Changes in transcription factor activity induce chromatin modifications, chromatin remodelling, a differential recruitment of the basal transcriptional machinery[1]. Modelling the dynamics of gene regulation is therefore essential to better understand why a cellular dynamic processes progresses through several steps, and what goes wrong in the case of disease.

The dynamics of gene regulation has classically been studied using time series data [2]. When dynamic processes progress asynchronously, such as in hematopoiesis, time series data are usually obtained by sorting different transition states and assessing bulk gene expression and transcription factor binding within the population [3, 4, 5, 6]. Alternatively, time series data can also be generated by synchronizing the dynamic process between cells. However, issues with time-resolution, heterogeneity and good in vivo synchronization models can often limit the predictive power of the dynamic models of gene regulation which can be constructed [2].

One of the main advantages of single-cell transcriptomics is the ability to quantify the exact cellular state of thousands of cells per experiment. The intercellular heterogeneity caused by naturally occurring biological stochasticity [7] can be exploited to predict regulatory interactions between transcription factors (TFs) and their target genes. The computational tools that infer gene regulatory networks (GRNs) from omics datasets are called network inference (NI) methods.

Several studies have highlighted how some regulatory interactions can be very dynamic while others show evidence of being static during consecutive developmental stages [8, 9]. Since regulatory interactions are context-dependent [10], attempting to create an accurate model of those processes by inferring a static regulatory network may have limited relevance. Case-wise NI methods¹ avoid predicting a static GRN and instead infer one GRN per cell (or per sample, for bulk omics data).

In order to compute a case-wise GRN for a single sample, Kuijjer et al. [11] and Liu et al. [12] employ similar strategies, namely by computing the difference of computing a static GRN for all the cases, and computing a static GRN for all the cases minus one. Since this procedure needs to be repeated for every case in the dataset, and because NI methods are already amongst the most computationally intensive analyses to perform on omics data, this methodology is not applicable for large omics datasets. Another case-wise NI method, SCENIC [13] infers case-wise GRNs by first inferring a static GRN using GENIE3 [14]. GENIE3 is a static NI method which uses Random Forests [15] feature importance scores to prioritise candidate regulators for a particular target gene. SCENIC then post-processes the static GRN to determine whether an interaction is enriched for particular cases, resulting in a case-wise GRN. In short, while several case-wise NI methods thus already exist, their implementation consisted of post-processing a static GRN to arrive at a case-wise GRN.

In this work, we introduce **bred**, the first ‘true’ case-wise NI method. For each interac-

¹Case-wise NI is sometimes also called sample-specific NI or case-specific NI.

tion in the inferred case-wise GRN, `brad` predicts both the regulatory strength and its effect. The case-wise GRNs – or ‘case-wise regulomes’ – can be analysed analogously to transcriptomics data; for example by clustering samples, inferring trajectories, or finding differentially activated regions. We demonstrate `brad` by applying it to a single-cell dataset of 22’122 hematopoietic cells from the Tabula Muris project [16], and to a collection of 14’963 bulk omics samples from The Cancer Genome Atlas project [17].

2 Results

2.1 Hematopoietic cells from Tabula Muris

From the Tabula Muris [16] project, we selected all cells involved in some part of the hematopoietic lineage tree. We computed case-wise GRNs between 22’122 remaining single cells, for the 2000 most variable target genes and a subset of transcription factors as regulators.

To summarise the similarity in GRNs across cases, we first visualise a dimensionality reduction of the case-wise GRNs, where every dot represents a single GRN (Figure 1A). The dimensionality reduction was computed by applying Fruchterman-Reingold [18] on the k -nearest-neighbour (k NN) graph of the case-wise GRN vectors. On the same k NN graph, the cells were clustered using the Louvain [19] clustering algorithm, and the clusters were labelled using prior information from Tabula Muris. For each cluster, we retained the 50 interactions with the highest mean importance scores (Figure 1B). In the visualisation, each node represents a gene, each edge an interaction, and the colour represents which cluster this interaction belongs to.

For the most part, clusters that are proximate in the dimensionality reduction (Figure 1A) are also closely connected in the GRN network view (Figure 1B). Notably, interactions from the different B cell clusters are almost not connected amongst each-other, while in the dimensionality reduction they are very proximate. Retaining more edges could result in the different B cells being connected in the GRN network view.

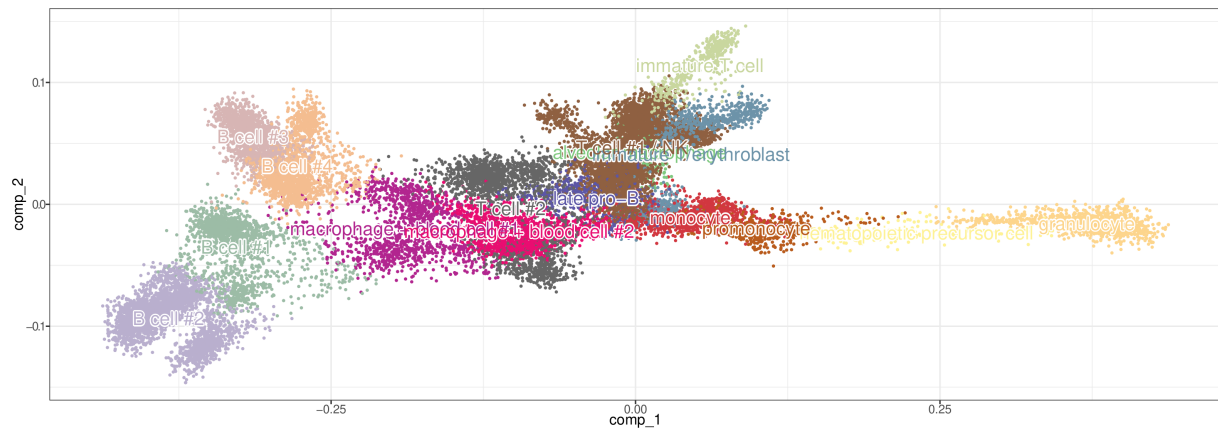
At the centre of the network lie many strongly connected interactions which are not specific for any particular cluster, but instead are common to almost all of the clusters. These are likely housekeeping genes or genes involved in key hematopoietic processes.

2.2 The Cancer Genome Atlas

We included all available RNA-seq profiles from The Cancer Genome Atlas. In total, we computed case-wise GRNs for 14’963 tumour samples from 50 sub-projects, including 44 different cancer entities. Clusters of case-wise GRNs were relabelled by the project they originate from, or by the cell types or organs the samples originate from.

Most groups of interactions are highly specific to just one cancer type. Samples from the CPTAC project were split into two groups, clustering together with LUAD (lung adenocarcinoma) and kidney carcinoma samples. According to the meta-information of CPTAC

A



B

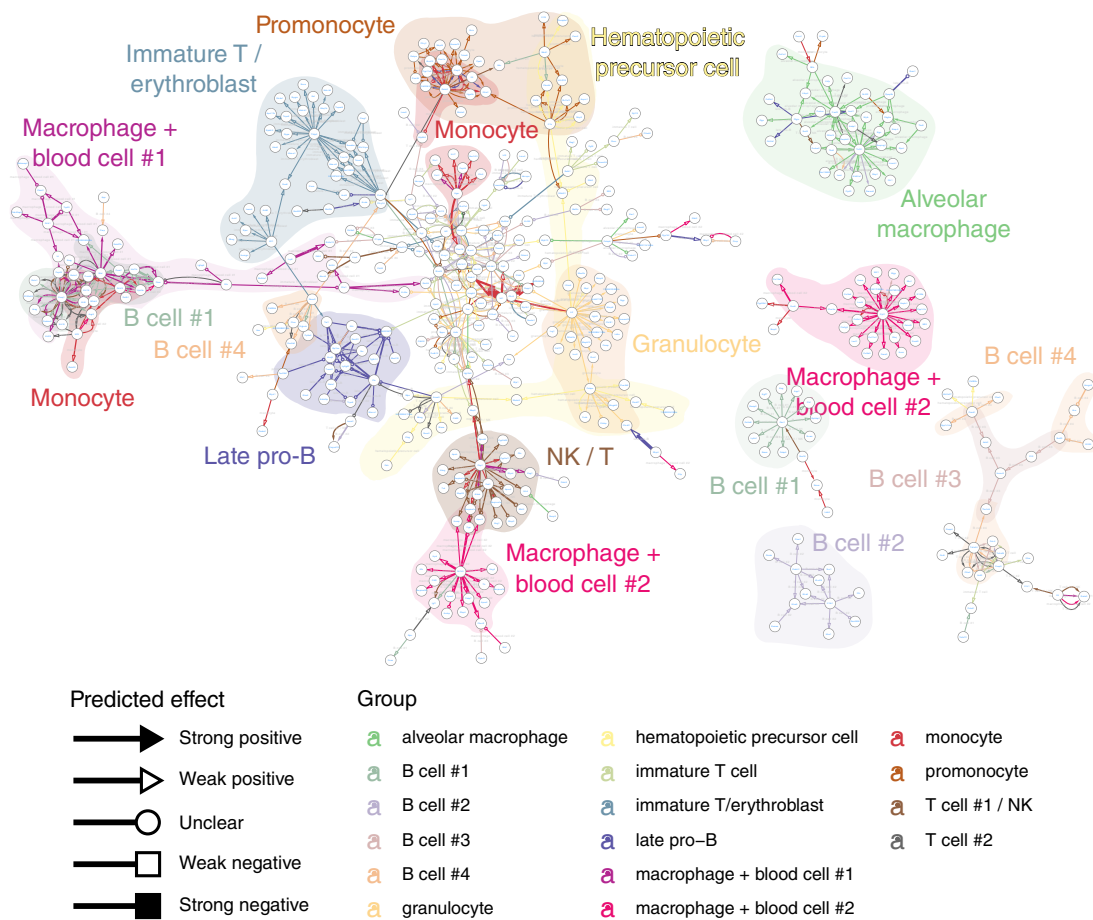
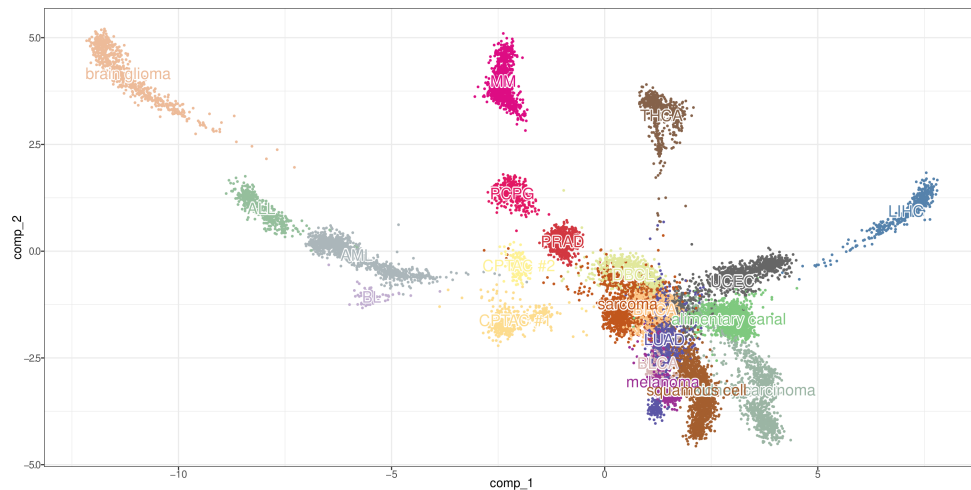


Figure 1: 22'122 case-wise GRNs of hematopoietic cells from the Tabula Muris project. A: Each dot in this dimensionality reduction corresponds to the GRN of a single cell. The dimensionality of the cell-specific regulome matrix was reduced using Fruchterman-Reingold and were clustered using Louvain clustering, after which the clusters were relabelled using the meta information from Tabula Muris. **B:** Per cluster, the 50 interactions with the highest mean importance score are visualised.

profiles, these samples consist of lung, kidney, and uterus adenocarcinomas, explaining the split in CPTAC samples.

A



B

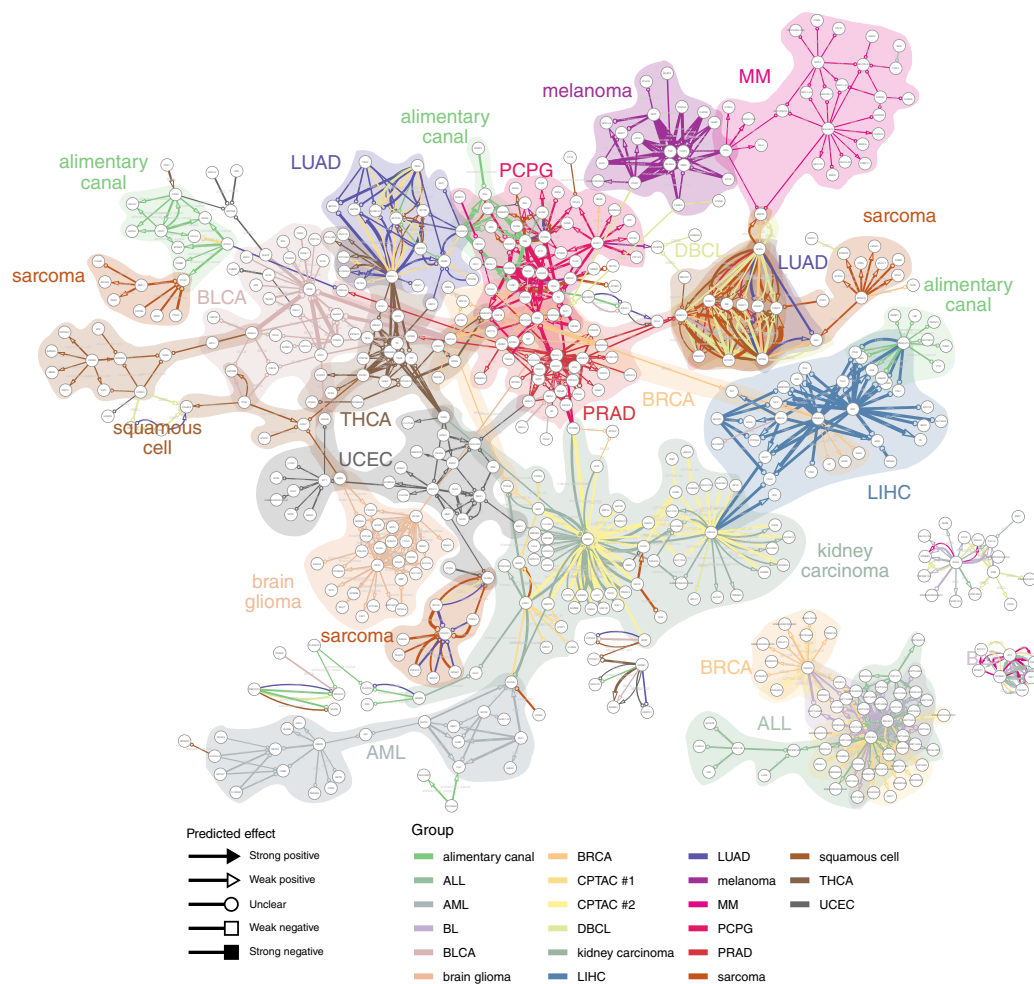


Figure 2: 14'963 case-wise GRNs of cancer bulk profiles from The Cancer Genome Atlas project. A: The different profiles are visualised and clustered according to their regulomic similarities. **B:** Visualisation of the strongest interactions per cluster shows both pathways distinct to particular cancer entities as well as pathways common to multiple cancer entities.

3 Discussion

The *bred* algorithm is a novel approach for directly computing case-wise GRNs for single-cell omics and bulk omics profiles alike. We applied the method to 22'122 hematopoietic cells and 14'963 cancer profiles, showing that interactions are often grouped together in modules specific for certain cell types, tissue types, or cancer types.

In this work, we only applied clustering methods to the case-specific 'regulome profiles', but other types of computational methods can be used to annotate and explore case-specific GRNs, such as trajectory inference and differential expression.

Going forward, we will provide further functional validation of the results generated by *bred*, as well as benchmark the algorithm against various real and *in silico* datasets.

4 Methods

4.1 Inferring case-wise GRNs

The task of inferring a static GRN (Figure 3A) can be reduced to a simpler problem, namely: for every target T , predict which of the potential regulators R_i regulate T (Figure 3B). This simplification allowed GENIE3 [14] to use Random Forest's [15] feature importance scores for inferring GRNs. Namely, a Random Forest is trained to predict the expression of a target gene of interest from the expression of potential regulators. The resulting Random Forest inherently allows to extract a feature importance score by observing the effect of each regulator in making a good prediction for the target expression. As in GENIE3, the target expression is first scaled to normalise feature importance scores across different targets.

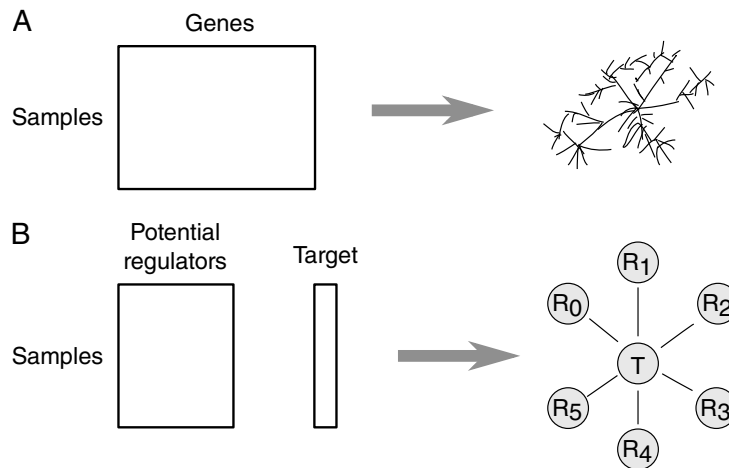


Figure 3: A: Inferring a gene regulatory network from an omics dataset can be reduced to a simpler problem. **B:** Given the expression of a target of interest and a set of potential regulators, predict which regulators regulate the target gene.

We make the same simplification in order to build case-wise GRNs, also using Random Forests to compute the feature importance scores. A Random Forest consists of K trees, each of which produces feature importance scores, and the feature importance scores of a forest is simply the mean feature importance scores of each of the trees.

Computing the case-wise feature importances of a tree consists of the following 8 steps (Figure 4). The ‘randomness’ of a Random Forest is due to only using a subset of the samples in the dataset in order to build a single decision tree. The samples are split into two groups, the ‘in-bag’ data and the ‘out-of-bag’ data (Figure 4A). A decision tree [20] is trained on the in-bag expression of the potential regulators in trying to predict the in-bag target gene expression (Figure 4B). The target expression of the out-of-bag samples is predicted using the decision tree (Figure 4C), and the squared error between the real and target expression is computed (Figure 4D). For each sample in the out-of-bag set, this vector represents how well the decision tree was able to predict the expression of the target gene.

The next few steps are repeated for every potential regulator R_i . Within the out-of-bag samples, the expression of R_i is randomly shuffled. The target expression of the out-of-bag samples is again calculated (Figure 4F), as well as the squared error between the real target expression and the predicted expression is calculated (Figure 4G). The importance of regulator R_i for an out-of-bag sample S_j is defined as the increase in squared error between the predicted target expression and the real target expression, after perturbing the expression of R_i (Figure 4H).

Steps F-G are repeated for every potential regulator R_i . By aggregating all of the feature importance scores over all the samples, regulators and targets, we obtain an M -by- N -by- P tensor².

A moderately-sized dataset could contain $M = 10'000$ samples, $N = 2'000$ regulators, and $P = 10'000$ target genes. Due to memory constraints, only interactions with an average importance value (across all samples) higher than a minimum threshold are retained.

To compute the case-wise GRNs, we implemented the abovementioned methodology in C++ in a modified version of the `ranger` R/C++ package [21].

4.2 Predicting the effect of an interaction

To predict the effect of a potential regulator R_i on a target gene T for a given tree, the Pearson correlation is calculated between the difference in regulator expression (before and after shuffling the values), and the difference in target expression prediction.

$$\begin{aligned} \text{effect}(R_i \rightarrow T) &= \text{cor}(x, y), \\ \text{with } x &= \text{expr_shuffled[:, } R_i] - \text{expr[:, } R_i], \\ \text{and } y &= \text{predict(tree, expr_shuffled)} - \text{predict(tree, expr)}. \end{aligned}$$

The Pearson correlation between two variables x and y is usually defined as shown in Equation 1. Computing r_{xy} for each (regulator, target) pairs, across all trees, would require storing large amounts of data.

²This is the origin of the name of the method, “bred”.

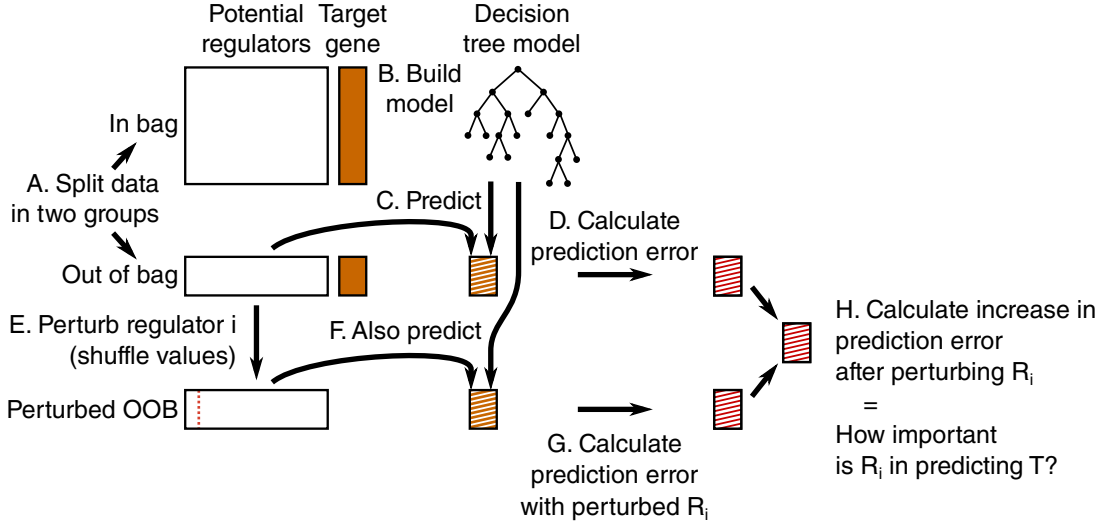


Figure 4: Calculating the feature importance score for one decision tree and one target consists of 8 distinct steps. A: Randomly split the data into two groups, the in-bag data and the out-of-bag data. **B:** The in-bag data is used to train a decision tree to try to predict the expression of the target gene from the expression values of the regulators. **C:** The decision tree is used to predict the gene expression of the target gene of the out-of-bag samples. **D:** Sample-specific squared error values are computed. **E:** Repeat steps E-H for every regulator R_i . Perturb the expression of regulator R_i in the out-of-bag samples. **F:** Again predict the gene expression of the target gene with the perturbed expression values. **G:** Again compute the sample-specific squared error values. **H:** The difference between the prediction error on the perturbed dataset versus the prediction error on the unperturbed is the importance in R_i in predicting T

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \times \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (1)$$

However, by rearranging the formula, it can be defined as Equation 2.

$$r_{xy} = \frac{\sum (x_i \times y_i) - \sum x \times \sum y / n}{\sqrt{(\sum x_i^2 - (\sum x)^2 / n)} \times \sqrt{(\sum y_i^2 - (\sum y)^2 / n)}} \quad (2)$$

For every regulator R_i during a perturbation in a given tree, only 6 values need to be stored, namely $A = \sum x_i$, $B = \sum y_i$, $C = n$, $D = \sum x_i \times y_i$, $E = \sum x_i \times x_i$, and $F = \sum y_i \times y_i$.

For every (regulator, target) pair, these values are summed, and the r_{xy} is calculated as shown in Equation 3.

$$r_{xy} = \frac{D - A \times B / C}{\sqrt{(E - A^2 / C)} \times \sqrt{(F - B^2 / C)}} \quad (3)$$

The following cutoffs were used to determine the effect.

- Strong negative: $r_{xy} < -0.4$
- Weak negative: $-0.4 \leq r_{xy} < -0.2$
- Unclear: $-0.2 \leq r_{xy} \leq 0.2$
- Weak positive: $0.2 < r_{xy} \leq 0.4$
- Strong positive: $0.4 < r_{xy}$

4.3 Clustering of case-wise GRNs

To perform downstream analysis on the cases, first a k -nearest neighbour (k NN) graph of the cases is computed. In order for the k NN graph to better emphasise similarities in GRNs rather than absolute euclidean distances, we first reduce the dimensionality of the case-by-interaction matrix to case-by-20 matrix using Landmark Multi-Dimensional Scaling [22] with a Spearman rank distance metric.

Next, KD-trees are used to calculate the k NN graph efficiently. The cases in the dataset are visualised and clustered using the Fruchterman-Reingold [18] and Louvain clustering [19], respectively.

The following R packages provided implementations for each of these algorithms: lmds, RANN, igraph [23].

4.4 Visualising clustered GRNs

After Louvain clustering, the interactions of the 50 interactions with highest mean importance per cluster are retained. These interactions are visualised in Cytoscape [24], in which nodes depict genes, edges depict predicted regulatory interactions, coloured according to which cluster they are predicted for. The shape of the arrow denotes the predicted effect of the regulatory interaction.

5 References

- [1] Antoine Coulon et al. "Eukaryotic Transcriptional Dynamics: From Single Molecules to Cell Populations". In: *Nature Reviews Genetics* 14 (July 9, 2013), p. 572. DOI: 10.1038/nrg3484.
- [2] Ziv Bar-Joseph, Anthony Gitter, and Itamar Simon. "Studying and Modelling Dynamic Biological Processes Using Time-Series Gene Expression Data". In: *Nat. Rev. Genet.* 13.8 (Aug. 2012), pp. 552–564.
- [3] Noa Novershtern et al. "Densely Interconnected Transcriptional Circuits Control Cell States in Human Hematopoiesis". In: *Cell* 144.2 (2011), pp. 296–309.
- [4] Gillian May et al. "Dynamic Analysis of Gene Expression and Genome-Wide Transcription Factor Binding during Lineage Specification of Multipotent Progenitors". In: *Cell Stem Cell* 13.6 (2013), pp. 754–768.
- [5] Vladimir Jojic et al. "Identification of Transcriptional Regulators in the Mouse Immune System". In: *Nat. Immunol.* 14.6 (2013), pp. 633–643. DOI: 10.1038/ni.2587. Identification.
- [6] Debbie K Goode et al. "Dynamic Gene Regulatory Networks Drive Hematopoietic Specification and Differentiation". In: *Dev. Cell* 36.5 (2016), pp. 572–587.
- [7] Olivia Padovan-Merhar and Arjun Raj. "Using Variability in Gene Expression as a Tool for Studying Gene Regulation". In: *Wiley Interdisciplinary Reviews. Systems Biology and Medicine* 5.6 (Nov. 2013), pp. 751–759. ISSN: 1939-005X. DOI: 10.1002/wsbm.1243. pmid: 23996796.
- [8] Victoria Moignard et al. "Characterization of Transcriptional Networks in Blood Stem and Progenitor Cells Using High-Throughput Single-Cell Gene Expression Analysis". In: *Nat. Cell Biol.* 15.4 (Apr. 2013), pp. 363–372.
- [9] Cristina Pina et al. "Single-Cell Network Analysis Identifies DDIT3 as a Nodal Lineage Regulator in Hematopoiesis." In: *Cell reports* 11.10 (2015), pp. 1503–1510. ISSN: 2211-1247. DOI: 10.1016/j.celrep.2015.05.016. pmid: 26051941.
- [10] Balázs Papp and Stephen Oliver. "Genome-Wide Analysis of the Context-Dependence of Regulatory Networks". In: *Genome Biology* 6.2 (Jan. 27, 2005), p. 206. ISSN: 1474-760X. DOI: 10.1186/gb-2005-6-2-206.
- [11] Marieke Lydia Kuijjer et al. "Estimating Sample-Specific Regulatory Networks". In: *iScience* 14 (Mar. 28, 2019), pp. 226–240. ISSN: 2589-0042. DOI: 10.1016/j.isci.2019.03.021. pmid: 30981959.
- [12] Xiaoping Liu et al. "Personalized Characterization of Diseases Using Sample-Specific Networks". In: *Nucleic Acids Research* 44.22 (2016), e164–e164. ISSN: 0305-1048. DOI: 10.1093/nar/gkw772. pmid: 27596597.
- [13] Sara Aibar et al. "SCENIC: Single-Cell Regulatory Network Inference and Clustering". In: *Nature Methods* (Oct. 2017). ISSN: 1548-7091. DOI: 10.1038/nmeth.4463.

- [14] Vân Anh Huynh-Thu et al. "Inferring Regulatory Networks from Expression Data Using Tree-Based Methods". In: *PLoS ONE* 5.9 (Jan. 2010), e12776. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0012776. pmid: 20927193.
- [15] Leo Breiman. "Random Forests". In: *Machine Learning* 45 (2001), pp. 5–32.
- [16] Nicholas Schaum et al. "Single-Cell Transcriptomics of 20 Mouse Organs Creates a Tabula Muris". In: *Nature* 562.7727 (Oct. 2018), pp. 367–372. ISSN: 1476-4687. DOI: 10.1038/s41586-018-0590-4.
- [17] John N Weinstein et al. "The Cancer Genome Atlas Pan-Cancer Analysis Project." In: *Nature genetics* 45.10 (Oct. 2013), pp. 1113–20. ISSN: 1546-1718. DOI: 10.1038/ng.2764. pmid: 24071849.
- [18] Thomas M. J. Fruchterman and Edward M. Reingold. "Graph Drawing by Force-Directed Placement". In: *Software: Practice and Experience* 21.11 (1991), pp. 1129–1164. ISSN: 1097-024X. DOI: 10.1002/spe.4380211102.
- [19] Vincent D Blondel et al. "Fast Unfolding of Communities in Large Networks". In: *Journal of Statistical Mechanics: Theory and Experiment* 2008.10 (Oct. 9, 2008), P10008. ISSN: 1742-5468. DOI: 10.1088/1742-5468/2008/10/p10008.
- [20] L Breiman et al. *Classification and Regression Trees*. Wadsworth Publishing Company, 1984.
- [21] Marvin N Wright and Andreas Ziegler. "Ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R". In: *Journal of Statistical Software* 77.1 (Mar. 2017). DOI: 10.18637/jss.v077.i01.
- [22] Seunghak Lee and Seungjin Choi. "Landmark MDS Ensemble". In: *Pattern Recognition* 42.9 (Sept. 2009), pp. 2045–2053. ISSN: 00313203. DOI: 10.1016/j.patcog.2008.11.039.
- [23] Gabor Csardi and Tamas Nepusz. "The Igraph Software Package for Complex Network Research". In: *InterJournal, Complex Systems* 1695.5 (2006), pp. 1–9.
- [24] Paul Shannon et al. "Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks". In: *Genome Research* 13.11 (Nov. 1, 2003), pp. 2498–2504. DOI: 10.1101/gr.1239303.