

# ANANSE: an enhancer network-based computational approach for predicting key transcription factors in cell fate determination

Quan Xu<sup>1</sup>, Georgios Georgiou<sup>1</sup>, Siebren Frölich<sup>1</sup>, Maarten van der Sande<sup>1</sup>,  
Gert Jan C. Veenstra<sup>1</sup>, Huiqing Zhou<sup>1,2,\*</sup> and Simon J. van Heeringen<sup>1,\*</sup>

<sup>1</sup>Radboud University, Department of Molecular Developmental Biology, Faculty of Science, Radboud Institute for Molecular Life Sciences, 6525GA Nijmegen, The Netherlands and <sup>2</sup>Radboud University Medical Center, Department of Human Genetics, Radboud Institute for Molecular Life Sciences, 6525GA Nijmegen, The Netherlands

Received November 04, 2020; Revised June 02, 2021; Editorial Decision June 23, 2021; Accepted June 28, 2021

## ABSTRACT

Proper cell fate determination is largely orchestrated by complex gene regulatory networks centered around transcription factors. However, experimental elucidation of key transcription factors that drive cellular identity is currently often intractable. Here, we present ANANSE (ANalysis Algorithm for Networks Specified by Enhancers), a network-based method that exploits enhancer-encoded regulatory information to identify the key transcription factors in cell fate determination. As cell type-specific transcription factors predominantly bind to enhancers, we use regulatory networks based on enhancer properties to prioritize transcription factors. First, we predict genome-wide binding profiles of transcription factors in various cell types using enhancer activity and transcription factor binding motifs. Subsequently, applying these inferred binding profiles, we construct cell type-specific gene regulatory networks, and then predict key transcription factors controlling cell fate transitions using differential networks between cell types. This method outperforms existing approaches in correctly predicting major transcription factors previously identified to be sufficient for *trans*-differentiation. Finally, we apply ANANSE to define an atlas of key transcription factors in 18 normal human tissues. In conclusion, we present a ready-to-implement computational tool for efficient prediction of transcription factors in cell fate determination and to study transcription factor-mediated regulatory mechanisms. ANANSE is

freely available at <https://github.com/vanheeringen-lab/ANANSE>.

## INTRODUCTION

Every multicellular organism develops from a single cell. During this process, cells undergo division and differentiation, eventually forming a diversity of cell types that are organized into organs and tissues. How one cell develops into different cell types, a process known as cell fate determination, is critical during development. It has been shown that transcription factors (TFs) play key roles in cell fate determination (1–6). TFs bind to specific cis-regulatory sequences in the genome, including enhancers and promoters, and regulate expression of their target genes (7,8). The interactions between TFs and their downstream target genes form gene regulatory networks (GRNs), controlling a dynamic cascade of cellular information processing which shapes cell fate and identity (9,10). Cell fate determination is orchestrated by a series of TF regulatory events, largely by complex GRNs (11). The key role of TFs and GRNs in cell fate determination is further corroborated by examples of cell fate conversions, often referred to as cellular reprogramming (12,13). Cellular reprogramming includes generating induced pluripotent stem cells (iPSCs) from somatic cells, and *trans*-differentiation that converts one mature somatic cell type to another without undergoing an intermediate pluripotent state (1–6). These reprogramming processes are initiated by enforced expression of combinations of key TFs, which is believed to alter the GRNs at the level of gene expression and the epigenetic landscape (14–16).

In the past, identification of key TFs driving cellular differentiation or reprogramming was often performed by experimental screening or testing candidate genes, which is labor-intensive and inefficient. Therefore, there is a need for better predictions of key TFs in cell fate determination,

\*To whom correspondence should be addressed. Tel: +31 24 3616850; Email: [s.vanheeringen@science.ru.nl](mailto:s.vanheeringen@science.ru.nl)  
Correspondence may also be addressed to Huiqing Zhou. Email: [j.zhou@science.ru.nl](mailto:j.zhou@science.ru.nl)  
Present address: Georgios Georgiou, Cergentis BV, Utrecht, the Netherlands.

which can help to understand developmental processes and serve to instruct experimental cellular reprogramming approaches. Different computational methods for predicting key transcription factors or master regulators in the context of cellular transitions have been reported. Some are based on gene expression differences between cell types (17–20). Other methods use GRNs in combination with expression differences to identify candidate key TFs (21–25). However, these GRNs are usually inferred based on a measure of co-expression (26,27), which requires many different samples and which cannot easily distinguish directionality.

The use of (predicted) transcription factor binding sites allows for directionality and has been shown to improve GRN inference (28–30). Mogrify is an example of a method that uses not only gene expression but also GRNs constructed based on TF binding motifs in promoters to predict TFs that are capable of inducing conversions between cell types (25). However, most GRN-based approaches that incorporate TF motifs only include promoters or promoter-proximal regulatory elements. It has been well established that TFs that control tissue- and cell type-specific gene expression in cell fate determination and development often bind to enhancers (15,31–33). Binding of tissue- and cell type-specific TFs largely to enhancers is also confirmed by a large number of genome-wide chromatin immunoprecipitation followed by sequencing analyses (ChIP-seq) (34,35), e.g. TP63 in keratinocytes and ZIC2 in embryonic stem cells (15,36). Furthermore, analysis of enhancers and enhancer clusters allows for identification of master regulators, corroborating their relevance in cell type-specific gene regulation (37,38). Large compendia of transcription factor binding profiles and enhancer-associated histone modifications can also be used to prioritize transcriptional regulators (39,40). Therefore, we reasoned that a computational method that uses enhancer properties to infer enhancer-based GRNs would improve the prediction of directed regulatory interactions at a genome-wide scale. Furthermore, most current computational tools require comprehensive training or background data, such as cell/tissue expression data or pre-constructed networks. This means that they cannot be applied in new biological contexts or in non-model species that are less well-studied. Finally, these datasets and the computational algorithms are not always publicly accessible, which prevents the general usage of these methods in studying transcriptional regulation or designing new *trans*-differentiation strategies.

Here, we established an enhancer GRN-based method, ANalysis Algorithm for Networks Specified by Enhancers (ANANSE), that infers genome-wide regulatory programs and identifies key TFs for cell fate determination. We predicted cell type-specific TF binding profiles with a model that incorporates activities and sequence features of enhancers. Second, combining TF binding profiles and gene expression data, we built cell type-specific enhancer GRNs in each cell type or tissue. We used reference GRNs, constructed from known TF-target gene interactions and experimental data of TF perturbations, to evaluate the quality of the inferred GRNs. Third, we predicted the key TFs underlying cell fate conversions based on a differential network analysis. Compared with other reported prediction algorithms, ANANSE recovers the largest fraction of TFs

that were validated by experimental *trans*-differentiation approaches. The results demonstrate that ANANSE can prioritize TFs that drive cellular fate changes. Finally, to demonstrate the wide utility of ANANSE, we applied it to 18 human tissues and generated an atlas of key TFs underlying human tissue identity.

## MATERIALS AND METHODS

### Analysis of the genomic distribution of TF binding sites

For every transcription factor, we combined all the peaks in the ReMap database (41) by taking the peaks in all cell types and tissues for this specific TF. TFs that had <600 peaks were removed. This resulted in a data set of ChIP-seq peaks from 296 unique TFs. The percentage of peaks in each genomic location was calculated using the ChIPseeker R package (version 1.20.0) (42). The fgsea R package (version 1.10.1) was used to do the gene set enrichment analysis (GSEA) (43).

We used the classification of the Human Protein Atlas (44) to determine tissue-specific TFs. This classification is divided in several groups based on gene expression patterns using RNA-seq from human tissues. We took the union of tissue enriched genes (at least a 5-fold higher FPKM level in one tissue compared to all other tissues), group enriched genes (5-fold higher average FPKM value in a group of 2–7 tissues compared to all other tissues) and tissue enhanced genes (at least a 5-fold higher FPKM level in one tissue compared to the average value of all 32 tissues).

### Datasets

All H3K27ac ChIP-seq, ATAC-seq, and RNA-seq data used in this study was obtained from GEO or ENCODE (45–63). For all data sets with ENCODE identifiers we downloaded the BAM files (ATAC-seq; H3K27ac ChIP-seq) or the FASTQ files (RNA-seq) from the ENCODE portal (64) (<https://www.encodeproject.org/>). For data sets with a GSM accession, FASTQ files were downloaded and further processed using seq2science (version v0.4.3) (65), see paragraph below. All data sets and accession numbers are summarized in Supplementary Table S1.

### ChIP-seq, ATAC-seq and RNA-seq analyses

Analysis of publicly available ChIP-seq, ATAC-seq and RNA-seq analysis was performed with seq2science (version v0.4.3) (65). Genome assembly hg38 was downloaded from UCSC with genomepy 0.9.1 (66). The reads of the ChIP-seq experiments were mapped to the human genome (hg38) using STAR (version 2.5.3a) with default settings (67). Duplicate reads were marked and removed using Picard (68). Peaks were called on the ChIP-seq data with only the uniquely mapped reads using MACS2 (version 2.7) relative to the Input track using the standard settings and a *q*-value of 0.01 (69). The measurement of consistent peaks between replicates was identified by IDR (version 2.0.3) (70).

ATAC-seq reads were trimmed with fastp (version 0.20.1) (71) and aligned with bwa-mem (version 0.7.17) (72) to the hg38 genome. Mapped reads were removed if they did not have a minimum mapping quality of 30, were a (secondary)

multimapper or aligned inside the ENCODE blacklist (73). Reads were shifted for tn5 bias. Duplicate reads were removed with picard MarkDuplicates (version 2.23.8) (68). Peaks were called with macs2 (version 2.2.7) (69) with options ‘-shift -100 -extsize 200 -nomodel -keep-dup 1 -buffer-size 10000’ in BAM mode.

Quantification of expression levels was performed on RNA-seq data, using salmon (version 0.13.0) (74) with default settings and Ensembl transcript sequences (version GRCh38 release-103) (75). Salmon’s transcript-level quantifications results were imported and aggregated to gene level counts by the tximport R package (version 1.12.3) (76). The expression level (transcript-per-million, TPM) of each cell type and the differential expression fold change between two cell types were calculated using the DESeq2 R package (version 1.24.0) (77). The expression TPM data used to predict key TFs for *trans*-differentiation is shown in Supplementary Table S2 and differential gene expression data is shown in Supplementary Table S3.

### Defining putative enhancer regions

To generate a collection of putative enhancer regions, we collected all transcription factor ChIP-seq peaks from ReMap 2018 ([http://remap.univ-amu.fr/storage/remap2018/hg38/MACS/remap2018\\_all\\_mac2\\_hg38\\_v1.2.bed.gz](http://remap.univ-amu.fr/storage/remap2018/hg38/MACS/remap2018_all_mac2_hg38_v1.2.bed.gz)) (41). We took the summit of all peaks and extended these 25 bp up- and downstream. Based on this file, we generated a coverage bedGraph using bedtools genomecov (78). We performed peak calling on this bedGraph file using bdgpeakcall from MACS2 (version v2.7.1) (69), with the following settings:  $l = 50$  and  $g = 10$ . We performed the peak calling twice, setting  $c$  to 4 and 30, respectively. All peaks from  $c = 30$  were combined with all peaks of  $c = 4$  that did not overlap with the peaks of  $c = 30$ . We then removed all regions on chrM and extended the summit of the peaks 100 bp up- and downstream to generate a final collection of 1 268 775 putative enhancers of 200 bp. This collection of enhancers is available at Zenodo with doi 10.5281/zenodo.4066423.

The coverage.table script from GimmeMotifs (version 0.15.3) (79,80) was used to determine the ATAC-seq and H3K27ac intensity, as expressed by the number of reads, in all enhancer peaks (2000 bp centered at the enhancer summit for H3K27ac; 200 bp for ATAC-seq). All counts were quantile normalized using qnorm (version 0.4.0) (81).

### Prediction of transcription factor binding

To train the ANANSE models we used ChIP-seq peaks for 237 TFs from REMAP in 6 cell types: hESC, HepG2, HeLa-S3, K562, MCF-7 and GM12878. ATAC-seq and H3K27ac ChIP-seq data for these cell types was downloaded from public repositories, see Supplementary Table S1. For both assays, the number of reads was determined in regions of 200 bp (ATAC-seq) or 2kb (H3K27ac) centered at the enhancer summit. Read counts were log-transformed and quantile normalized. To test the prediction performance of the ANANSE model a cross-validation procedure was used. For each TF, models were trained on binding in all enhancers, except those on chromosomes chr1,

chr8 and chr21 (held-out chromosomes). The evaluation was only performed on those TFs for which peaks in multiple cell types were available. Each cell type was left out (held-out cell types) and the classifier was trained on data of the other cell type(s). In this manner, performance metrics (ROC AUC and PR AUC) were calculated based on enhancers located on held-out chromosomes in held-out cell types.

Binding was predicted using four type(s) of input features: TF motif scores, ATAC-seq signal in enhancers, H3K27ac ChIP-seq signal in enhancers and (optionally) the average ChIP-seq signal of REMAP peaks in enhancers. ANANSE uses a standard logistic regression model as implemented in scikit-learn (82). Equation (1) shows an example of a model, using all four types of input.

$$\log \frac{p_{f,l}}{1 - p_{f,l}} = \beta_1 S_{f,l} + \beta_2 E_{ATAC,l} + \beta_3 E_{H3K27ac,l} + \beta_4 E_{ChIP,l} \quad (1)$$

where  $p_{f,l}$  is the probability of a transcription factor  $f$  binding to enhancer  $l$ .  $S_{f,l}$  is the highest motif z-score of all motifs associated with transcription factor  $f$  in enhancer  $l$  and  $E_{ATAC,l}$ ,  $E_{H3K27ac,l}$  and  $E_{ChIP,l}$  represent the enhancer intensity of enhancer  $l$ , based on scaled and normalized ATAC-seq signal, scaled and normalized H3K27ac ChIP-seq signal and average REMAP ChIP-seq signal, respectively.

ANANSE incorporates a flexible selection of models, the choice of which depends on the type of input that is available. The minimal input consists of the motif score and either ATAC-seq or H3K27ac ChIP-seq signal in enhancers.

The non-redundant database of 1796 motifs was created by clustering all vertebrate motifs from the CIS-BP database using GimmeMotifs (80,83) as described in (79). The GimmeMotifs package (version 0.15.3) (79,80) was used to scan for motifs in enhancer regions. The GC normalization setting in GimmeMotifs package was used to normalize the GC% bias in different enhancers. To correct for the bias of motif length, z-score normalization was performed on the motif scores. Normalization was done per motif, based on motif matches in random genomic regions using the same motif scan settings. The highest z-score was chosen if a TF had more than one motif.

### Gene regulatory network inference

The weighted sum of the TF binding probability, predicted on the basis of the enhancer intensity and the motif score, within 100kb around TSS is defined as the TF–gene binding score (Equation 2). The distance weight is based on a linear genomic distance between the enhancer and the TSS of a gene according to Equation (3).

$$B_{x,r} = \sum_k w_{k,r} s_{k,x} \quad (2)$$

where  $B_{x,r}$  is the binding score between TF  $x$  and target gene  $r$ ,  $w_k$  is the weighted distance between an enhancer and the target gene and where  $s_k$  is predicted binding intensity at genomic position  $k$  of TF  $x$ . The distance weight calculation

was similar to the method previously described in (84), except that only the signal in predetermined enhancers is used and the weight of enhancers within 5kb of the TSS is set to 1.

$$w_k = \begin{cases} 1, & k \in (0 \text{ kb}, 5 \text{ kb}] \\ \frac{2e^{-\mu|k-t_r|}}{1+e^{-\mu|k-t_r|}}, & k \in (5 \text{ kb}, 100 \text{ kb}] \end{cases} \quad (3)$$

where  $t_r$  is the genomic position of the TSS of gene  $r$  and the parameter  $\mu$ , which determines the decay rate as a function of distance from the TSS, is set such that an enhancer 10 kb from the TSS contributes one-half of an enhancer within 5 kb from TSS.

We determined a measure of genome-wide TF activity,  $A_x$ , based upon the motif activity. The motif activity for all TF motifs was calculated based on ridge regression as implemented in scikit-learn (82) using GimmeMotifs 0.15.3 (79). Here, the motifs scores were used as input to predict either ATAC-seq and/or H3K27ac ChIP-seq signal. The TF activity is the maximum activity of the motifs associated with a TF, where the motif activity is defined as the mean of the ATAC-seq motif coefficients and the H3K27ac ChIP-seq coefficients.

The expression level of the TF  $E_x$  and the target gene  $E_r$ , expressed as transcripts per million (TPM), and the TF activity  $A_x$  and TF-gene binding score  $B_{x,r}$  were ranked and scaled, from 0 to 1, where 0 represents the lowest value and 1 represents the highest value. For ranking the TF expression, only the expression levels of TFs were used. The interaction score was calculated (Equation 4) by mean averaging the individual ranked scores (mean rank aggregation).

$$I_{x,r} = \frac{1}{4}(F(E_x) + F(E_r) + F(B_{x,r}) + F(A_x)) \quad (4)$$

where  $I_{x,r}$  is the interaction score between TF  $x$  and target gene  $r$  and  $F(X)$  represents the rank aggregated and scaled score. Ideally, the contributions of these individual scores would be determined by a supervised method, such as a linear regression, however, due to the lack of a high-quality gold standard reference data set we chose to combine the scores through mean averaging. To create the network, ANANSE uses *dask* (<https://dask.org>) (85) and *pyranges* (86).

### Gene regulatory network evaluation

We obtained GRNs from different sources. GRNBoost2, as implemented in *arboreto* (version 0.1.5), was used with default settings to infer networks from GTEx data. The GTEx expression data (GTEx\_Analysis\_2017-06-05\_v8\_RNASEQCv1.1.9\_gene\_tpm.gct.gz) was downloaded from the GTEx portal (<https://www.gtexportal.org/home/datasets>). Tissue-specific PANDA networks were downloaded from <https://sites.google.com/a/channing.harvard.edu/kimberlyglass/tools/gtex-networks>. Tissue-specific networks inferred from single cell data using SCENIC (87) were downloaded from <http://www.grndb.com/> (88). GRNs inferred from GTEx data with *corto* and ARACNE were downloaded from <https://giorgilab.org/corto-the-correlation-tool/> (26,29,89).

To evaluate the quality of the predicted GRNs, four different types of reference datasets were used: gene co-

expression, Gene Ontology (GO) annotation (The Gene Ontology, 2019), four regulatory interaction databases (DoRothEA (90), RegNetwork (91), TRUST (92) and MSigDB C3 (93)) and differential expression measurements after TF perturbations. The expression correlation database was downloaded from COXPRESdb (94), and the original mutual rank correlation score was scaled to 0 to 1 for each TF, with 1 being the highest and 0 the lowest, and all scaled correlation score higher than 0.6 or 0.8 were considered as true interaction pairs. The human GO validation Gene Association File (GAF) (version 2.1) was downloaded from <http://geneontology.org>. We used all TF-gene pairs that were annotated with at least one common GO term as true positives. The TF perturbation data set was obtained by downloading the 'TF\_Perturbations\_Followed\_by\_Expression' data set from Enrichr (95,96). For the random network we used the same network interaction structure, but here we randomized the interaction score (the edge weight) by permutation of the scores. The AUC of ROC and PR for each cell type GRN and corresponding random GRN were calculated.

### Influence score inference

To calculate the influence score for the transition from a source cell type to a target cell type, we used the GRNs for both cell types. In each network, we selected the top 100k interactions based on the rank of its interaction score. We obtained a differential GRN by taking the interactions only located in the target cell type. The difference of the interaction score was used as the edge weight for the differential GRN.

Based upon the differential GRN a local network was built for each TF, up to a maximal number of three edges. Using Equation (5), a target score was calculated for each node in the network, based on 1) its edge distance from the TF of interest, 2) the interaction score and 3) the change in expression between the source cell type and the target cell type.

$$N_x^s = \sum_{r \in V_t} |G_r^s| \frac{P_{x,r}^s}{L_{x,r}^s} \quad (5)$$

where  $r \in V_t$  is each gene ( $r$ ) in the set of nodes ( $V_t$ ) that make up the local sub-network of TF  $x$ . In other words,  $V_t$  represents all target genes that are directly or indirectly targeted by TF  $x$ . To incorporate indirect target genes, only genes up to three steps away are considered. This distance (number of edges or steps) is represented by  $L_{x,r}^s$ , the level (or the number of steps) that gene  $r$  is away from TF  $x$  in the network  $s$ . Nodes located further from the TF have less effect on the target score.  $P_{x,r}^s$  is the interaction score between TF  $x$  and target gene  $r$  and  $G_r^s$ , the expression score, is the log-transformed fold change of the expression of gene  $r$ .

The target score ( $N_x^s$ ) for each TF is the sum of the scores from all the nodes in its local network. Nodes present in multiple edges are calculated only for the edge closest to the TF of interest. Self-regulating nodes are not considered. The target score and the  $G_r^s$  of each TF are scaled to 0 to 1, and the mean of them was defined as the influence score of

this TF. Subsequently, all TFs are ranked by their influence score.

### Trans-differentiation evaluation

To evaluate the performance the ANANSE influence score calculation we used key TFs from *trans*-differentiation experiments. We compared ANANSE results to previously reported methods: Mogrify, LISA, BART, VIPER, CellNet and the method of D'Alessio *et al.* (20,21,24,25,39,40). Mogrify and Mogrify full prediction results were downloaded from <https://moglify.net/>. For LISA (version 1.2), all differentially expressed genes from fibroblast to each target cell type were used as input (39). For BART, we uploaded the top 1000 differentially expressed genes to <http://bartweb.org/>, using BART 2.0. The DoRothEA network was downloaded from <https://github.com/saezlab/dorothea>. All differentially expressed genes from fibroblast to each target cell type and networks (ANANSE or DoRothEA) were used as input of VIPER (version 1.24.0). The CellNet predictions were obtained from (25). The prediction results of the method of D'Alessio *et al.* were obtained from the original paper (20). Results were compared with the experimentally validated TFs as true positives and all other TFs as false positives.

### Regulatory profile analysis of human tissues

The RNA-seq data of 18 human tissues was downloaded from (<https://www.proteinatlas.org/humanproteome/tissue>) (97). The H3K27ac ChIP-seq and ATAC-seq accession numbers are listed in Supplementary Table S1. The gene expression score of each tissue was calculated by taking the log<sub>2</sub> TPM fold change between a tissue and the average of all other tissues. The GRN of each tissue was inferred using ANANSE. For prediction for TFs of one tissue, GRN interaction scores of all other tissues were averaged as the source GRN. All correlation analyses were clustered using hierarchical clustering. The modular visualization of anatograms and tissues was done using the gganatogram package (version 1.1.1) (98).

## RESULTS

### Cell type-specific transcription factors predominantly bind to enhancers

To systematically examine TF binding patterns in the genome in relation to cell type specificity, we downloaded the binding sites of 296 human TFs from the ReMap project, which re-analyzed all publicly available ChIP-seq data in various cell types and tissues (41). To determine the genomic distribution of these binding sites, we divided the genome into different genomic categories according to human UCSC known gene annotation (99), and assigned binding sites to these categories based on the locations of the binding sites (Figure 1). We grouped these categories into two main classes: (i) a promoter-proximal class, containing promoter ( $\leq 2$  kb), 5' UTR and 1st exon peaks, and (ii) a promoter-distal class, referred to as 'Enhancers', containing all exons except the first, the first intron, other introns and intergenic categories. The percentage of TF

binding sites in each genomic category was calculated, and TFs were ordered according to their percentages in the promoter-proximal class (Figure 1A) (Supplementary Table S4).

As expected, we found that the majority of TFs (77.5%) mainly bind in *cis*-regulatory regions that are distal from the promoter (Figure 1A). These binding sites will not necessarily all be functional, however, they are not close to gene promoters and contain the majority of the enhancers. For the purpose of this study we will refer to them as enhancers. However, different TFs show different binding distributions, with a preference in either the promoter range or in the enhancer range (Figure 1B). Given the relevance of enhancers in cell type-specific gene regulation, we reasoned that cell type-specific TFs would have a larger fraction of peaks in enhancers than constitutively expressed TFs and performed Gene Set Enrichment Analysis (GSEA) (43) on TF expression in different tissues. We defined tissue-specific TFs using previously established categories based on gene expression patterns using RNA-seq from human tissues, including tissue-enriched genes, group-enriched genes, and tissue-enhanced genes (Human Protein Atlas; see methods for details) (44) (Figure 1C). GSEA showed that TFs mostly binding to enhancers are enriched for tissue-specific expression (adjusted  $P$  value =  $2.0e-4$ ) (Figure 1C) (Supplementary Table S4). For example, SOX10 is a critical TF during neural crest and peripheral nervous system development (100), while TP63 is a master regulator in epithelial development (58). Both of these tissue-specific TFs showed a very high percentage of enhancer-binding, 93% for SOX10 and 82% for TP63 (Figure 1A).

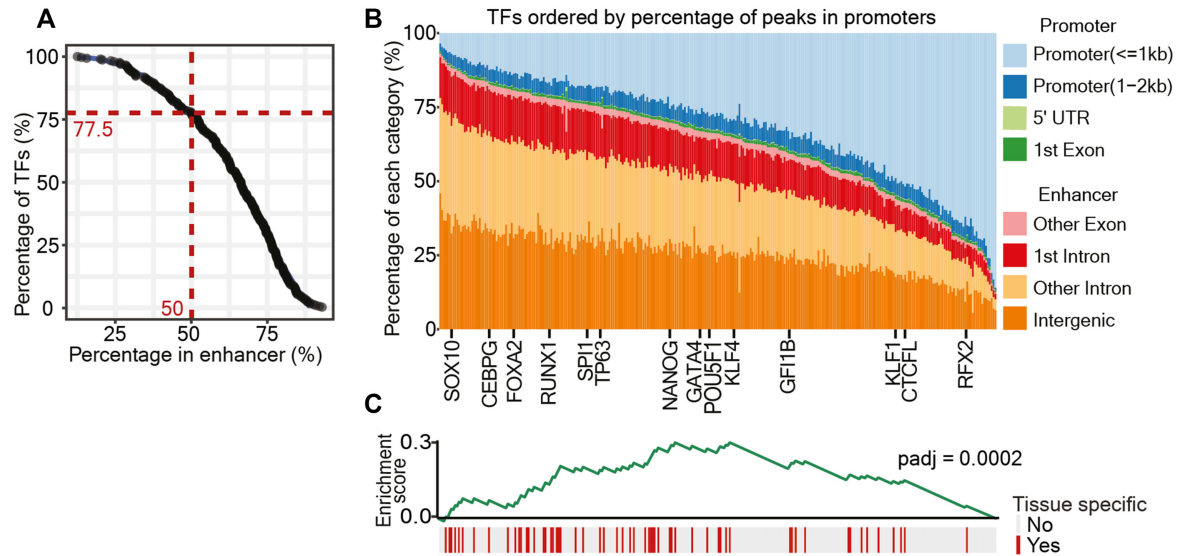
Taken together, our analysis of transcription binding sites confirmed that distal *cis*-regulatory elements are especially relevant for tissue-specific TFs. This emphasizes that including enhancer information in computational methods for predicting key TFs in cell fate determination could be highly beneficial.

### ANANSE: an enhancer network-based method to identify transcription factors in cell fate changes

Starting from the premise that the majority of TFs predominantly bind to enhancer regions, we developed ANANSE, a network-based method that uses properties of enhancers and their GRNs to predict key TFs in cell fate determination (Figure 2). As *trans*-differentiation is an ideal model for studying cell fate conversions controlled by key TFs, we set out to use this model to validate our computational approach. In the following paragraphs a conceptual overview of ANANSE is provided. Subsequently we will validate each of the steps involved.

First, we inferred cell type-specific TF binding profiles for each cell type (Figure 2A). The input data of ANANSE consists of genome-wide measurements of enhancer activity (defined below) and transcription factor motifs. We inferred the TF binding probability based on a supervised model that integrates the enhancer activity combined with TF motif scores.

Second, we constructed cell type-specific GRNs based on the inferred TF binding probability, the transcription factor activity, and the expression levels of the TF and pre-



**Figure 1.** Tissue-specific TFs predominantly bind to enhancers. (A) The percentage of TF binding sites in putative enhancers. The human genome was split into several categories: Promoter ( $\leq 1$  kb), Promoter (1–2 kb), 5' UTR and first exon, other exons, first intron, other introns and intergenic; these categories were further grouped into a promoter-proximal class (Promoter ( $\leq 1$  kb), Promoter (1–2 kb), 5' UTR and first exon) and an enhancer class (other exons, first intron, other introns and intergenic). Out of 296 human TFs, 77.5% have at least 50% of their binding sites in the enhancer class of the genome. (B) Genomic location analysis of binding sites of 296 human TFs. The percentage of binding sites of each TF in different categories (as described in A) was calculated, and indicated with different colors. TFs were ordered by the percentage of binding sites within the promoter-proximal class. Several example TFs are marked at the bottom of the figure. (C) Gene Set Enrichment Analysis (GSEA) on tissue-specific TFs and their enhancer binding. The red bars mark the tissue-specific TFs. The order of TFs is consistent with (B). Gray bars represent TFs that do not show tissue-specific gene expression. The GSEA enrichment score is represented by the green line ( $P_{\text{adj}}: 2.0e-4$ ).

dicted target genes (Figure 2B and C). The nodes in the network represent the TFs and their target genes. In this network, a TF node can also be a target gene of another TF. The TF–gene interaction scores, represented by edges of the network, are calculated based on the predicted TF binding probability, the distance between the enhancer and the target gene, the genome-wide TF activity and the expression of both the TF and the target gene. By integrating these data, ANANSE determines the interaction score of each TF–gene pair.

Third, we calculated the ‘influence score’ (21,25), a measure of importance of a TF in explaining transcriptional differences between two cell types (Figure 2D). In this step, the difference in TF–gene interaction scores (the inferred networks, 2C) between the source and the target cell types is calculated. This differential network is combined with the expression differences between the cell types to determine the influence score.

The details of the algorithms are described in the following sections.

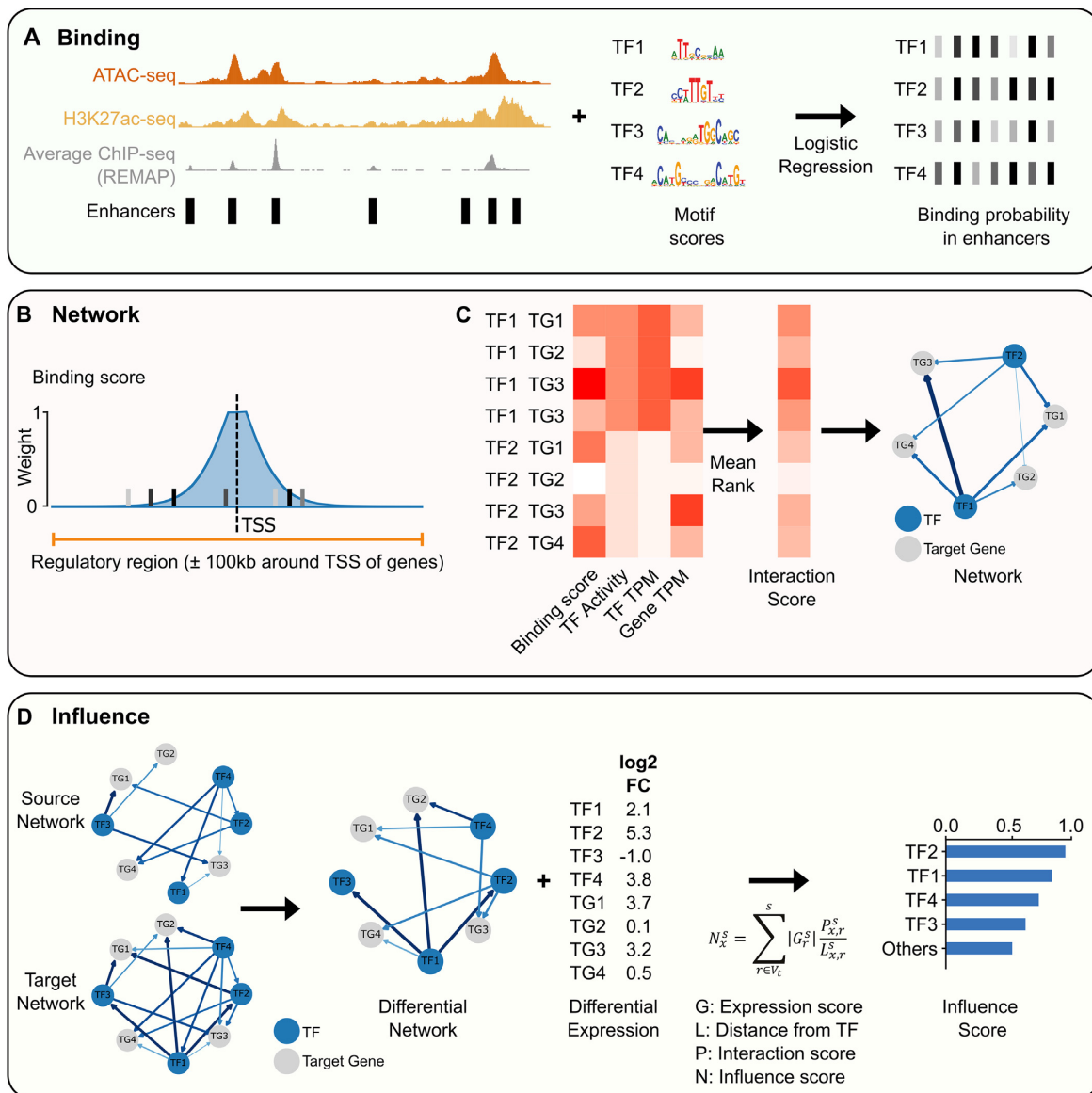
### Transcription factor binding can be predicted by the motif score in combination with the enhancer activity

Sequence-specific TFs bind to their cognate DNA motifs in the genome and activate or repress their target genes. To infer the target genes of a TF, the genomic binding sites of this TF are informative. ChIP-seq has been broadly used to identify TF binding sites at a genome-wide scale. However, it is unfeasible to perform ChIP-seq for every TF in all cell types, e.g. due to the availability and quality of the TF anti-

bodies. Therefore, it would be highly beneficial to be able to predict binding sites of individual TFs in a given cell type.

Here, we used a conceptually simple logistic regression classifier to predict the TF binding probability in putative enhancers based on the TF motif  $z$ -score, the enhancer activity and (optionally) the average TF binding signal in these regions (see Materials and Methods for details). Our model uses a predefined set of putative enhancers as input. In this work, we used a set of 1.3 million putative enhancer regions based on an integration of all TF ChIP-seq peaks from ReMap 2018 (41). Alternatively, putative enhancers can be based on genome-wide measurements that provide relatively accurate estimates of enhancer locations, such as ATAC-seq, DNaseI-seq or EP300 ChIP-seq. The enhancer activity is based on two genome-wide assays: chromatin accessibility as measured by Assay for Transposase Accessible Chromatin with high-throughput sequencing (ATAC-seq) (101) and the presence of the post-translational histone modification H3K27 acetylation as measured by ChIP-seq.

To train and evaluate our model, we used ChIP-seq peaks of 237 TFs in six cell lines (hESC, Hep-G2, HeLa-S3, K562, MCF-7 and GM12878) from REMAP (41). We examined the locations of the TF peaks by overlapping with our enhancer reference, and only the subset of peaks that overlapped with these enhancer regions was kept. We downloaded and mapped public ATAC-seq and H3K27ac ChIP-seq data for these cell lines (see Supplementary Table S1). For both assays, we determined the number of reads in regions of 200 bp (ATAC-seq) or 2kb (H3K27ac) centered at the enhancer summit. Read counts were log-transformed and quantile normalized. For each enhancer, we scanned for



**Figure 2.** An overview of the ANANSE method. ANANSE consists of three different modules: binding prediction, network inference and influence score calculation. **(A)** TF binding prediction using a supervised model. The binding probability of all TFs with an associated motif is calculated using on the basis of four input data types: 1) a set of reference cis-regulatory regions (here based on REMAP (41) ChIP-seq integration), 2) genome-wide enhancer activity measurements (ATAC-seq and/or H3K27ac ChIP-seq), 3) the REMAP average ChIP-seq intensity and 4) TF motif scores. The bars in the right panel represent the predicted binding probability of four TFs in the enhancers shown in the left panel. **(B)** A schematic overview of the first step in gene regulatory network inference, calculation of the TF–gene binding score. A binding score of each TF and gene combination is calculated by aggregation of all enhancers near a gene, weighted by a distance function. The orange line shows 100 kb up- and downstream of the TSS of the corresponding target gene, the range that is used to include enhancers for calculation. The bars represent the predicted TF binding probabilities within the 100 kb range around the gene. The height of the shaded light blue area represents the weight calculated based on the linear genomic distance from TSS of the target gene to the enhancers (84). For example, the distance weight is 1 for the distance of 1 kb from the TSS, and 0 for the distance of 100 kb from the TSS. **(C)** A schematic overview of the gene regulatory network inference using rank aggregation. The heatmap on the left represents the input for each TF and target gene (TG) combination: the binding score (according to B), the genome-wide TF activity, the TF expression level (transcripts-per-million; TPM) and the target gene expression level (TPM). All four scores are ranked and scaled from 0 to 1, and the mean of the four scores of each TF–gene pair is defined as the interaction score (right heatmap) of the corresponding TF–gene pair. **(D)** Overview of the influence score calculation. The influence score represents how well the expression differences between two cell types can be explained by a TF. First, a differential GRN is calculated between source and target cell type (left). Then, the influence score is calculated based on the gene expression log<sub>2</sub> fold change, the distance from the TF to the gene in the predicted network, and the interaction score in the differential network between TF and gene (middle). The barplot (right) shows the ranked influence score of all TFs calculated from the differential GRN.

motifs in a 200 bp region centered at the peak summit using GimmeMotifs (79,80). The motif z-score was calculated by GimmeMotifs with the GC%-normalization option. The log-odds score based on the positional frequency matrix is normalized by using the mean and standard deviation of scores of random genomic regions. These random regions are selected to have a similar GC% as the input sequence.

The goal of the binding model in ANANSE is to predict binding for all TFs, based on a supervised model. However, not all TFs have ChIP-seq available for training. Therefore, we implemented a two-pronged approach. We trained a TF-specific supervised model for each TF for which we had training data, but we also trained a general classifier based on all TFs. In this manner, we can use a more performant TF-specific model for TFs that have ChIP-seq training data, but can still predict binding for all other TFs, as long as they have an associated motif. In both cases, the input data for the final trained model is identical, however, the TF-specific models will be better tuned to the binding patterns of their associated TF. To test the prediction performance of our model, we established a stringent cross-validation procedure (102). For each TF we trained on binding in all enhancers, except those on chromosomes chr1, chr8 and chr21 (held-out chromosomes). In addition, we only included the TFs in the evaluation for which we had more than one cell type available. In turn, each cell type was left out (held-out cell types) and the classifier was trained on data of the other cell type(s). In this manner, performance metrics are calculated based on enhancers located on held-out chromosomes in held-out cell types. We evaluated the performance of the ANANSE binding model using the AUC (Area Under Curve) of the Precision Recall curve (PR) as well as the (Receiver Operating Characteristic) ROC curve (Figure 3A). For comparison, we included two baselines. The ‘random’ baseline represents the performance that would be observed by ‘random guessing’ (0.5 for the ROC AUC; the proportion of positives in the evaluation set for the PR AUC). The more stringent ‘Average ChIP-seq’ baseline represents the performance that would be observed if the binding is predicted only by the number of different TFs in REMAP that bind to an enhancer (i.e. the predicted binding probability is directly proportional to the number of REMAP peaks overlapping an enhancer). We compared two versions of the ANANSE model, one with the average ChIP-seq peak signal of REMAP included, and one where this is not included (Figure 3A; ‘With average’ and ‘Without average’, respectively). The model where the average signal is not included is more representative of the performance on other reference enhancer sets, or in other species. The median PR AUC of 0.28 is significantly higher than that of the random baseline (median PR AUC 0.02,  $P$  Wilcoxon  $< 1e-58$ ) as well as the average baseline (median PR AUC 0.18,  $P$  Wilcoxon  $< 1e-38$ ). When we include the average signal, the performance is significantly improved (median PR AUC 0.38,  $P$  Wilcoxon  $< 1e-39$ ). The full ANANSE model also improves on the individual components, as models based on motif scores (median PR AUC 0.06), ATAC-seq (median PR AUC 0.19) or H3K27ac alone (median PR AUC 0.13), while higher than the random baseline, do not significantly outperform the average ChIP-seq baseline. Finally, the model performs well, even when only one of the

assays is used (median PR AUC of 0.28 and 0.31 for ATAC-seq and H3K27ac respectively).

For comparison to methods that train more complex supervised transcription factor-specific models, we also validated our model on the validation chromosomes (chr1, chr8, and chr21) in the validation cell types of the ENCODE-DREAM transcription factor binding challenge (Available from: <https://www.synapse.org/ENCODE>) (103). As a comparison, we used the Virtual ChIP-seq predictions (104). This is a newly developed supervised artificial neural network method to predict individual TF binding, which shows comparable results compared to the top ENCODE-DREAM entries. In this evaluation (Supplementary Figure S1A and B), our model scores better for some factors, such as CEBPA in liver and MAX in liver and K562 cells, while Virtual ChIP-seq performs better for other factors, most notably CTCF. This illustrates that the binding prediction of ANANSE is comparable to state-of-the-art approaches. A caveat here is that other methods, such as Virtual ChIP-seq, predict binding genome-wide, while ANANSE needs a set of putative enhancers as input. An advantage of the relatively simple model of ANANSE is that it generalizes to other TFs and other species, which will not have the abundance of training data provided by ENCODE for mouse and human. As another evaluation, we compared the ANANSE binding predictions to predictions based on DNase I footprinting (Supplementary Figure S1C and D) (105). Compared to DNase I footprints, ANANSE predictions have higher recall at the same precision (median 0.78 vs 0.01).

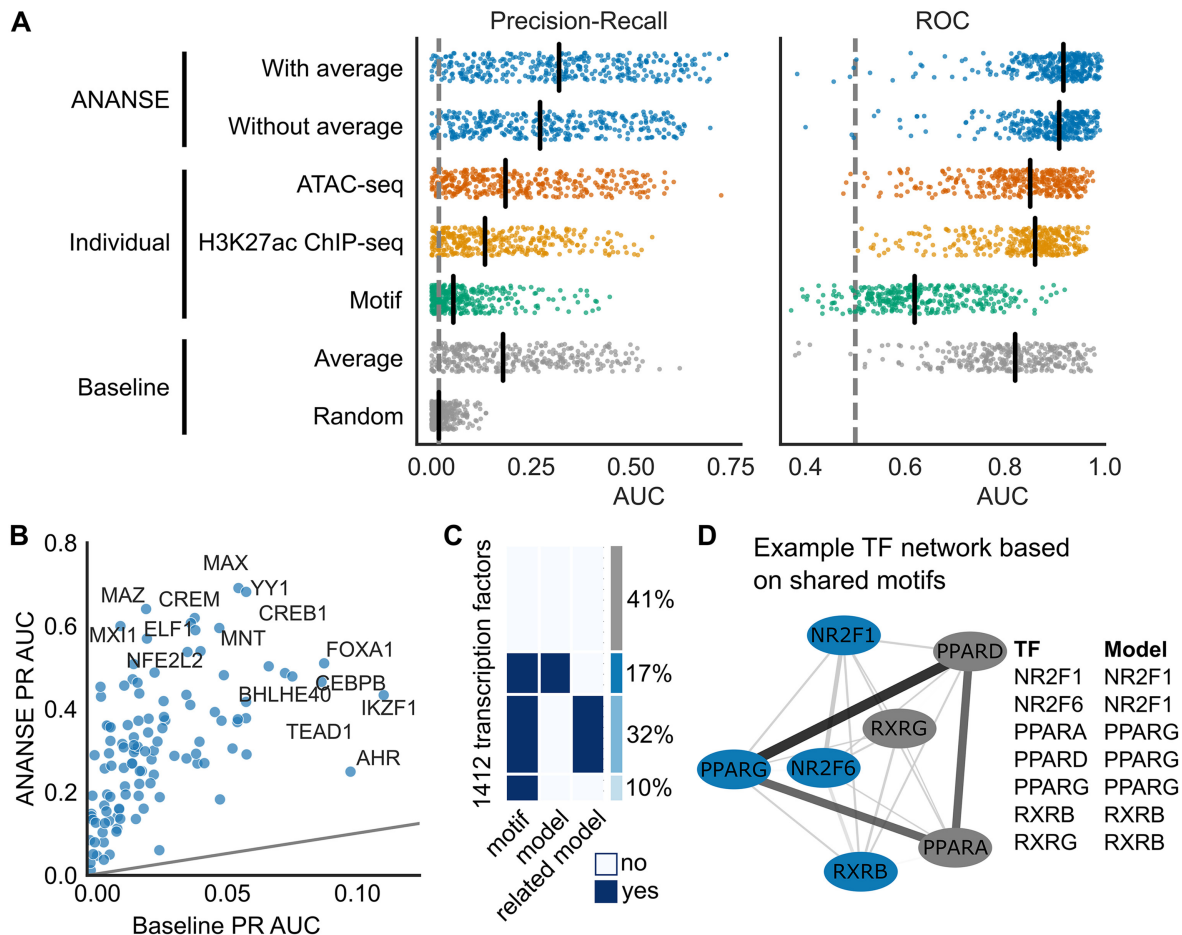
In total, these analyses illustrated that we established a TF binding site prediction method, which can be applied on the basis of one or two experimental measurements (ATAC-seq and/or H3K27ac ChIP-seq) and which shows state-of-the-art performance in prediction of TF binding.

### ANANSE predicts cell type-specific gene regulatory networks

Using the inferred cell type-specific binding profiles, we sought to determine the interactions of TFs and their target genes (TF-gene) to establish cell type-specific GRNs, represented by the TF-gene interaction score. To calculate these scores, we first identified all enhancers for each target gene and their associated binding scores. In our TF binding prediction model, we used H3K27ac ChIP-seq and ATAC-seq as training data. For each gene, we took all enhancers that are located within a maximum distance of 100 kb of the transcription start site (TSS). Subsequently, the strength of a TF-gene interaction in the network was defined by the sum of the predicted TF binding strength in all identified enhancers of the target gene weighted by the distance (Figure 2B), similar to the regulatory potential (84). The distance weight was calculated from the linear genomic distance between the enhancer and the TSS of a gene, such that distal enhancers receive a low weight and nearby enhancers have a high weight (Figure 2B). This model resulted in a TF-gene binding score, indicating the TF-target gene binding intensity for all combinations of TF and target gene pairs.

Next to the incorporation of the genome-wide measurements of enhancer activity via the TF-gene binding score,





**Figure 3.** The performance of predicting TF binding sites using TF motif scores and enhancer activities. **(A)** Evaluation of the TF binding prediction model in ANANSE. Shown is the area under the curve (AUC) of the Precision-Recall (PR; left) and Receiver-operator characteristic (ROC; right) of the prediction performance using REMAP ChIP-seq peaks as a reference. Plotted is the performance of 237 TFs in six cell lines. PR AUC and ROC AUC metrics were calculated using cross-validation. The performance is compared to two baselines (grey): random (proportion of positives for PR, 0.5 for ROC) and the average number of REMAP TF ChIP-seq peaks per enhancer. Performance on individual input data types is shown as reference: ATAC-seq (orange), H3K27ac ChIP-seq (yellow) and motif scores (green). Two ANANSE predictions (integration of ATAC-seq, H3K27ac ChIP-seq and motif scores; blue) are compared, based on the inclusion of the average REMAP ChIP-seq signal. Both models perform better than the baselines. **(B)** The scatterplot shows the improvement of the full ANANSE model compared to the random baseline, with some example factors highlighted. **(C)** An overview of how many human TFs have an associated model in ANANSE. Out of 1412 TFs, 17% have a trained model available, 32% of TFs have a model trained on a related TF (based on motif similarity, see D) and 10% have a motif and can use the general model. The remaining 41% of TFs do not have an associated motif. **(D)** An example of determining related TFs for model sharing. The network illustrates the similarity between a selection of nuclear receptors, as determined by the Jaccard index of their associated motifs (edge color and size). There is no trained model for PPARD, but it shares many motifs with PPARG, so it uses the PPARG model weights with PPARD motif scores.

we also calculated a measure of TF activity directly from the genome-wide enhancer activity. Here, we used a method similar to the motif activity response analysis (106,107) as implemented in GimmeMotifs (79). In this approach, the enhancer activity (ATAC-seq and H3K27ac signal) is modeled as a linear function of motif scores using penalized regression. The coefficients of the motif scores can be interpreted as an estimate of motif activity. For all TFs we used the maximum activity of all associated motifs as TF activity.

Finally, based on the assumption that the interaction of every TF–gene pair in a specific cell type is proportional to their relative expression, we included the expression level of the TF and the target gene, the TF and target expression scores. We ranked the expression level of the TF and the target gene, initially expressed as transcripts per million (TPM)

within each cell type to a normalized expression between 0 and 1, with the lowest expression as 0 and highest as 1.

To calculate the TF–gene interaction score, we combined the TF–gene binding score, the TF activity, and the TF and target expression scores using mean rank aggregation (Figure 2C). This score represents the strength of the regulatory interaction between a TF and a target gene. In this approach, a ‘target gene’ can also be a TF gene; a TF–gene interaction can represent a TF regulating the expression of a TF gene. Together, all TF–gene interaction scores represent a cell type-specific GRN.

To evaluate the quality of the GRNs inferred by ANANSE, we created GRNs for 15 tissues: adrenal gland, brain, colon, esophagus, heart, liver, lung, ovary, pancreas, prostate, skeletal muscle, skin, small intestine, spleen and

stomach. We collected ATAC-seq, H3K27ac and RNA-seq from public repositories (see Supplementary Table S1) and predicted binding profiles and GRNs using ANANSE. As comparison, we included five other GRN inference methods. We downloaded ARACNE, corto and PANDA networks created using GTEx expression data (26,29,89). We downloaded the GTEx expression data and created GRNs using GRNBoost2 (108). Finally, we downloaded GRNs created by SCENIC (87) with tissue single cell data as input from GRNdb (88).

To provide a comprehensive overview of GRN quality, we used four different types of reference datasets to calculate performance metrics: (i) regulatory interaction databases containing known TF–target gene interactions, (ii) differential expression measurements after TF perturbation, (iii) gene co-expression data and (iv) Gene Ontology (GO) annotation (109). We obtained the TF–gene interactions from four databases of regulatory interactions, DoRothEA (90), RegNetwork (91), TRRUST (92) and MSigDB C3 (93). DoRothEA is a gene set resource network containing different types of TF and target interactions. For this comparison, we only used the literature-curated interactions. RegNetwork is an integrated database of transcriptional and post-transcriptional regulatory networks in human and mouse, TRRUST is an expanded reference database of human and mouse transcriptional regulatory interactions and MSigDB C3 is a collection of gene sets that represent potential targets of regulation by TFs. While the TF–target gene databases are usually curated, they contain relatively few interactions. As a source of a more genome-wide validation we downloaded the ‘TF.Perturbations.Followed.by.Expression’ data set from Enrichr (95,96). This is a curated collection of genes that significantly change expression after TF perturbation. We downloaded co-expression data for human genes from the COXPRESdb database (94). All TF–gene pairs with either a correlation  $\geq 0.6$  or  $\geq 0.8$  were used as true positives. Finally, we used TF–gene pairs that were annotated with at least one common GO term as true positives.

To compare with the tissue-specific GRNs inferred by other methods, we selected per reference the Cartesian product of TFs and target genes of the sets of TFs and target genes in the reference. If a specific interaction was not present in the inferred GRN, we used the minimum interaction score. Supplementary Note S1 contains more details on the benchmark procedure. We evaluated the GRNs by calculating the PR AUC and ROC AUC, as compared to the reference interactions (Figure 4A and B; Supplementary Figure S2). For the ANANSE networks, the median AUC ranges from 0.61 using the RegNetwork reference (Supplementary Figure S2A) to 0.77 using DoRothEA (Figure 4A), while the median AUC of randomized networks is close to 0.5. When we compared ANANSE with other published GRN inference methods, all five methods show a significantly lower AUC using DoRothEA (Figure 4A) or TF perturbation references (Figure 4B). This holds true for all other references, except for the correlation reference, where PANDA scores higher (Supplementary Figure S2A). Some of the reference databases contain very few interactions (the positives in this evaluation) as compared to all possible in-

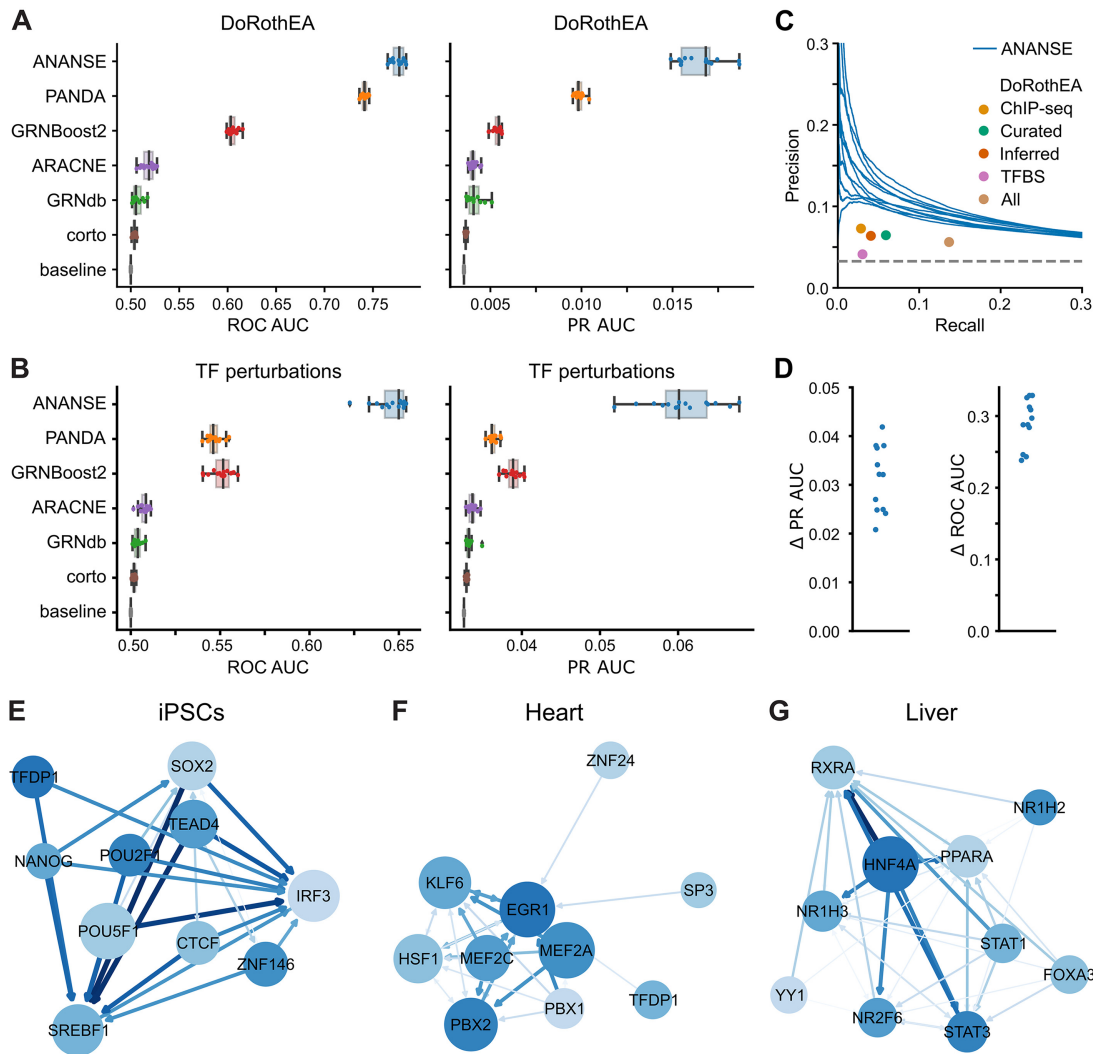
teractions (which determine the negatives). For instance, the fraction of positive interactions is 0.09% in TRRUST, 0.07% in RegNetwork, and 0.33% in MSigDB C3. Therefore, we also evaluated the predicted networks using the PR AUC (Figure 4D and Supplementary Figure S2B). In absolute terms, the PR AUC is considerably lower than the ROC AUC, especially for the DoRothEA, TF perturbation, MSigDB C3, TRRUST and RegNetwork reference sets (median PR AUC of 0.0168, 0.0601, 0.0297, 0.0123 and 0.0187, respectively), but for all tissues there is a relatively large and statistically significant difference between the predicted GRN and the randomized network ( $p$  Wilcoxon =  $9.86e-22$ ) (Figure 4A and B, Supplementary Figure S2B). The expression-based networks of the other five published methods show significant lower mean PR AUC when compared with ANANSE for all seven different types of reference datasets (Figure 4A and B and Supplementary Figure S2B).

To further evaluate the ANANSE networks, we compared them to all different types of interactions present in DoRothEA, using the TF perturbations as a reference. In Figure 4C the PR curves of the inferred networks for the different tissues are plotted, and the different interactions present in DoRothEA are represented by dots. These include interactions predicted by TF ChIP-seq binding near genes (ChIP-seq), curated interactions from the literature (Curated), interactions inferred using ARACNE-VIPER (Inferred) and interactions predicted by TF motifs scores in the promoter (TFBS). In addition the union of all sets is shown (All). The difference in recall at the same precision, and precision at the same recall between the ANANSE networks and the union set (All) of DoRothEA is shown in Figure 4D. These results illustrate that, using this benchmark, ANANSE-inferred networks better predict genes deregulated by TF perturbation.

The inclusion of enhancer information in the GRN inference is one of the unique features of ANANSE. To test if incorporating enhancer activity indeed leads to an improved performance, we compared the enhancer-based approach of GRN with an expression-based model and two promoter-based models (Supplementary Figure S3). Overall we found that the incorporation of enhancer information leads to better network inference, as tested by these benchmarks. The choice of distance and the method of combining enhancers has a smaller effect, with a 100 kb distance performing marginally better than 250kb and the distance-weighted sum performing better than the mean or sum with a uniform weight (Supplementary Figure S4).

To qualitatively assess the cell type-specific GRNs predicted by ANANSE, we chose one well-studied cell type (iPSC) and two tissues (heart and liver) and constructed their GRNs using the top ten predicted TFs of each cell type, as ranked by outdegree. The GRN of iPSCs includes well-known pluripotency factors such as POU5F1, NANOG and SOX2 (Figure 4E). The GRNs of heart and liver tissues contain marker genes, such as the myocyte factors MEF2A and MEF2C in heart (Figure 4F), and HNF4A and FOXA3 in liver (Figure 4G).

Taken together, our benchmarks and examples demonstrate that GRNs generated by ANANSE allow for meaningful cell type-specific prioritization of TFs.



**Figure 4.** Prediction of tissue-specific enhancer gene regulatory networks. (A) Evaluation of the predicted networks using the curated interactions of the DoRoThEA database. The left panel shows the AUC of ROC for 15 different tissues in a boxplot, with individual tissues marked as dots. The right panel shows the PR AUC. The ANANSE predicted networks (blue) are compared to other GRN inference approaches trained on GTEx expression data for the same tissues (PANDA in orange, GRNBoost2 in red; ARACNE in purple and corto in brown) and on the GRNdb networks inferred using SCENIC on single cell RNA-seq data from the same tissues (green). The random baseline is shown in gray. (B) The same evaluation as in (A) using a reference of differentially expressed genes after TF perturbation. (C) Comparison of the tissue-specific GRNs inferred by ANANSE to the different types of interactions in DoRoThEA, using the TF perturbations as a reference. The PR curves of the inferred networks for the different tissues are plotted (ANANSE; blue), and the different interactions present in DoRoThEA are represented by dots: interactions predicted by TF ChIP-seq binding near genes (orange), curated interactions from the literature (green), interactions inferred using ARACNE-VIPER (red) and interactions predicted by TF motif scores in the promoter (purple). The union of all DoRoThEA interactions is shown in brown. (D) The difference of the ANANSE GRNs with the union of DoRoThEA interactions in (C) expressed as the difference in precision at the same recall (left panel) and the difference in recall at the same precision (right panel). (E) Example network predicted for iPSCs. The blue circles show the top 10 TFs in this cell type, ranked by the outdegree in the top 100 000 edges. The size of the circle indicates the target gene number of the corresponding TF. The size and color of the blue arrows are relative to the interaction score between the two TFs. The color of the circle indicates the expression level of the corresponding TF. (F) Example network predicted for heart, visualized as in (E). (G) Example network predicted for liver, visualized as in (E).

### ANANSE accurately predicts key transcription factors for *trans*-differentiation

Having established that ANANSE-inferred GRNs can enrich for biologically relevant regulatory interactions, we aimed to use these GRNs to identify key TFs that regulate cell fate determination. To this end, *trans*-differentiation is a good model, as experimentally validated TFs have been determined for various *trans*-differentiation strategies. Here, we first inferred the GRNs for all cell types using our

ANANSE approach. The ANANSE-inferred GRN difference between a *source* and a *target* cell type, was calculated to represent the differential GRN between two cell types, which contains the GRN interactions that are specific for or higher in the target cell type. Subsequently, using an approach inspired by Mogrify (25), we calculated the influence score of TFs for these *trans*-differentiations by determining the differential expression score of its targets weighted by the regulatory distance (see Methods for details).

**Table 1.** The summary of seven experimentally validated *trans*-differentiations from fibroblast to target cell types. Experimentally validated TFs that were identified by ANANSE are highlighted in bold

Target cell type	Experimentally validated TFs	Reference
Astrocyte	NFIA, NFIB, <b>SOX9</b>	(110)
Cardiomyocyte	<b>GATA4</b> , <b>MEF2C</b> , TBX5	(111)
	HAND2, <b>NKX2-5</b> , <b>GATA4</b> , <b>MEF2C</b> , TBX5	(112)
Hepatocyte	ATF5, PROX1, FOXA2, <b>FOXA3</b> , <b>HNF4A</b>	(113)
	<b>FOXA1</b> , <b>FOXA3</b> , <b>HNF4A</b>	(114)
iPSC	<b>SOX2</b> , <b>POU5F1</b> , KLF4, MYC	(5)
	<b>POU5F1</b> , <b>SOX2</b> , NANOG, LIN28	(116)
	<b>POU5F1</b> , <b>SOX2</b>	(115)
Keratinocyte	<b>TP63</b> , <b>GRHL2</b> , TFAP2A, MYC	(117)
Macrophage	<b>CEBPA</b> , <b>CEBPB</b> , <b>SPI1</b>	(119)
	<b>SPI1</b> , <b>CEBPA</b> , <b>CEBPB</b>	(118)
Osteocyte	<b>RUNX2</b>	(120)
	<b>RUNX2</b> , <b>POU5F1</b> , MYCL	(121)

To evaluate the prediction by ANANSE, we used experimentally validated TFs for several *trans*-differentiation strategies. For this, we collected TFs for seven *trans*-differentiation strategies with fibroblasts as the source cell type. The target cell types include astrocytes (110), cardiomyocytes (111,112), hepatocytes (113,114), iPSCs (5,115,116), keratinocytes (117), macrophages (118,119), and osteocytes (120,121) (Table 1). We used ATAC-seq and H3K27ac ChIP-seq data of these cell types to create cell type-specific GRNs, and then calculated TF influence scores and ranked the TFs in each cell type.

When we calculate TF influence scores from cell type-specific GRNs, it is important to decide what size of GRN should be chosen in terms of the top number of edges. We inferred the key TFs for the seven *trans*-differentiations using six different sizes of GRNs (10K, 50K, 100K, 200K, 500K and 1M edges; Supplementary Figures S5 and S6, Supplementary Table S5). These results show that the approach is relatively invariant to the GRN size, with performance starting to decrease at 1 million edges. Here, we chose a GRN size of 100K interactions for all following analyses.

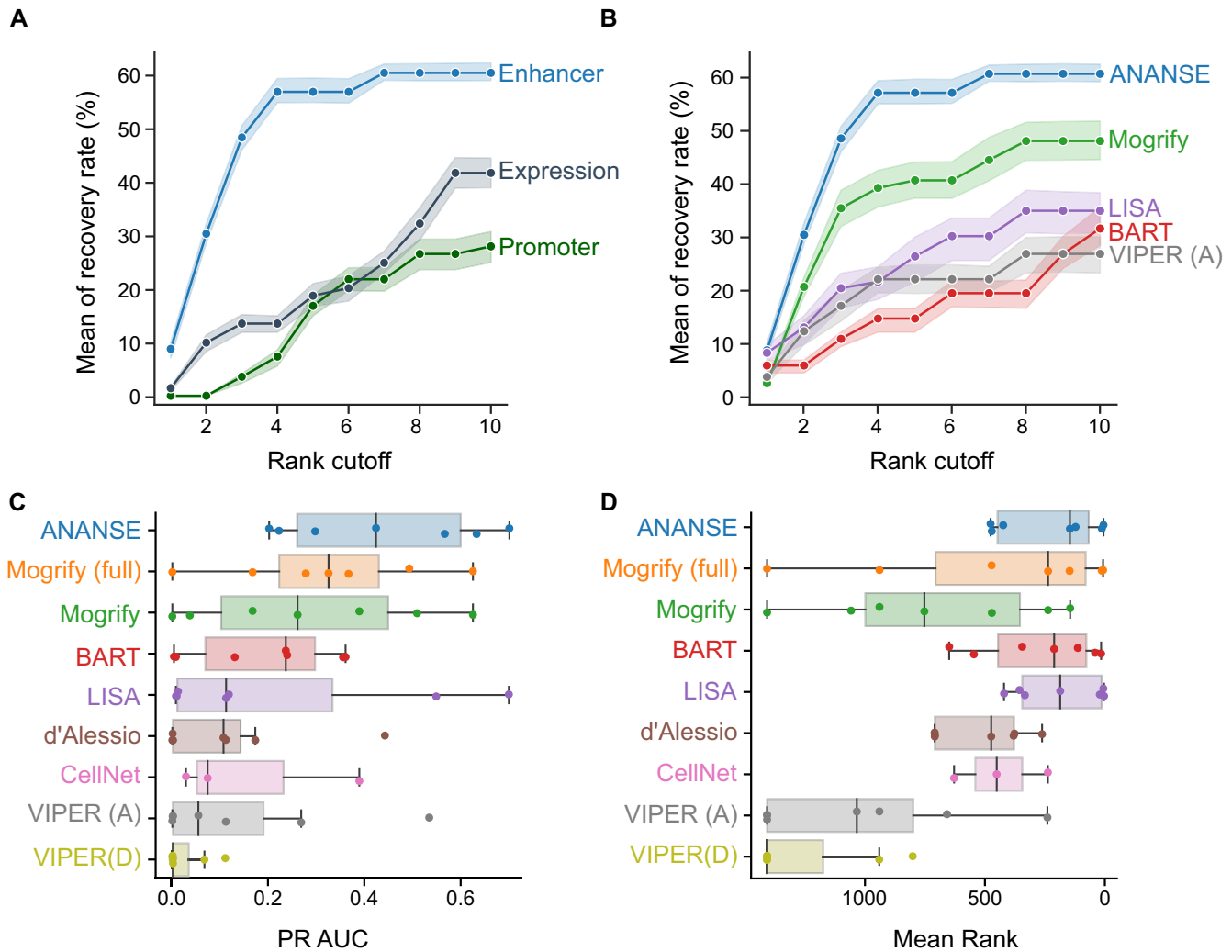
Using GRNs with 100K edges, we predicted the top 10 TFs for all seven *trans*-differentiations. In all cases, many of the experimentally defined TFs are included in the top 10 factors predicted by ANANSE (Table 1). For example, ANANSE predicts CEBPA, CEBPB and SPI1 for reprogramming fibroblasts to macrophages (118,119) and HNF4A, FOXA1 and FOXA3 for reprogramming to hepatocytes, which are consistent with the experimental *trans*-differentiation strategies (113,114).

To evaluate if the inclusion of enhancer information in the GRN inference in ANANSE resulted in more accurate predictions of TFs for *trans*-differentiation, we compared the enhancer-based approach of ANANSE with a promoter-based model, which does not take into account promoter-distal regulatory elements, and with a model based on only gene expression data (Figure 5A). We created both expression and promoter based GRNs of the seven source and target cell type combinations. For expression-based GRNs, we used only the mean of the scaled TPM of TFs and genes together as the interaction score of TFs and

genes. For the promoter-based GRNs, we selected the highest binding score of TFs within 2 kb of the TSS of the corresponding gene as the binding score of the TF–gene pair. Subsequently, the mean of the scaled TPM of the TF and the gene together with the binding score determines the interaction score of the TF and gene (Figure 4B). We then inferred the key TFs for the seven *trans*-differentiations using ANANSE and these two types of GRNs. The ANANSE influence score based on the enhancer GRNs includes 57% of the known TFs in the top four predictions (Figure 5A and Supplementary Table S6). In contrast, using the influence score based on the promoter GRN or the expression GRN, we could recover only 5% and 14% of the known TFs in the top four predictions (Figure 5A, Supplementary Figures S7 and S8A and Supplementary Table S6). These results demonstrate that using enhancers in the construction of GRNs significantly improves the prediction of relevant TFs in cell fate determination.

Next, we further quantified the performance difference between ANANSE and previously reported methods, namely Mogrify, LISA, BART, VIPER, CellNet and the method of D'Alessio *et al.* (20,21,24,25,39,40) (Figure 5B and Supplementary Figure S8B). For Mogrify, we downloaded both the prioritized list of TFs based on TF expression in source cell types and GRN overlap, as well as the full unfiltered list of TFs. For VIPER, we predicted TFs with the DoRothEA network, 'VIPER (D)', and with the GRNs inferred by ANANSE, 'VIPER (A)'. For these comparisons, we aimed to include all seven *trans*-differentiation strategies. In some cases, as data for the exact cell type is unavailable, similar cell or tissue types were used as surrogates. For example, the osteoblast-Sciencell was used to substitute for osteoblast. For CellNet, we used the previously described results of three cell types: hepatocytes, iPSCs and macrophages (25). For LISA, we used all differentially expressed genes from fibroblast to each target cell type as input (39). For all methods, the resulting TFs were ranked according to the relevant output score of the method. Using the seven cell type conversions as a reference, ANANSE has the highest recovery at all rank cutoffs up to 10 (Figure 5B, Supplementary Figure S8B and Supplementary Figure S9). ANANSE predicts a mean of 57% TFs using the top four TFs ranked by influence score, while other methods predict a maximum of 39% of TFs with this rank cutoff (Figure 5B, Supplementary Figure S8B and Supplementary Figure S9). When the number of predicted TFs was increased to ten, ANANSE could increase its recovery rate to 61%, while the maximum mean recovery of other methods is 47% (Figure 5B, Supplementary Figure S8B and Supplementary Figure S9). In addition to the mean recovery rate, we also evaluated the PR AUC (Figure 5C) and the mean rank of all known *trans*-differentiation factors (Figure 5D). In these analyses, ANANSE shows the highest median score (PR AUC or mean rank). However, other methods perform nearly as well in these benchmarks, such as Mogrify (both with PR AUC and mean rank) and BART and LISA (mean rank).

In summary, these analyses show that including enhancers in the GRN construction significantly improves the prediction of TFs in cell fate conversion and that ANANSE outperforms other established methods, based on experimentally validated *trans*-differentiation TFs. Our results



**Figure 5.** Evaluation of the performance of ANANSE using experimentally validated *trans*-differentiation strategies. (A) The line plots show the comparison of the predicted top TFs for *trans*-differentiation from cell type-specific networks. Based on the difference between two networks, TFs were prioritized using the influence score calculation implemented in ANANSE. Shown is the fraction of predicted TFs compared to all known TFs based on *trans*-differentiation protocols described in the literature (y-axis) as a function of the top number of TFs selected (x-axis). The mean of recovery rate is the average of all TF sets when the corresponding *trans*-differentiation has several different experimental validated TF sets. The shaded area represents the minimum and maximum percentage of corresponding recovered TFs when using six out of seven *trans*-differentiations. Three different types of networks were used: gene expression (dark blue), promoter-based TF binding in combination with expression (dark green), and enhancer-based TF binding in combination with expression (blue). (B) The line plots show the comparison of the predicted top TFs for *trans*-differentiation based on different computational methods. The y-axis indicates the percentage of experimentally validated cell TFs that are recovered as a function of the number of top predictions, similar as in A). Six different methods are shown: ANANSE (blue), Mogrify (green), LISA (purple), BART (red) and VIPER with ANANSE networks (gray). The shaded area represents the minimum and maximum percentage of corresponding recovered TFs when using six out of seven *trans*-differentiations. Mogrify and CellNet only contain the top 8 predicted factors. For visualization purposes, only a subset of the evaluated methods is shown. The remaining methods are shown in Supplementary Figure S8B. (C) The PR AUC of the same different *trans*-differentiations as in A and B, shown as a boxplot. Individual *trans*-differentiations are shown as dots. (D) The mean rank of the experimentally determined factors of the same *trans*-differentiations as in A and B, shown as a boxplot. Individual *trans*-differentiations are shown as dots.

demonstrate that ANANSE can prioritize biologically relevant TFs in cell fate determination.

#### ANANSE identified an atlas of key transcription factors in normal human tissues

The gene expression programs that drive the cellular differentiation programs of different tissues are largely controlled by TFs. To find out which key TFs drive cell fate determination in different tissues, we applied ANANSE to

a much wider range of human tissue data. We downloaded H3K27ac ChIP-seq data of 18 human tissues from the ENCODE project (46) and the RNA-seq data of corresponding tissues from the Human Protein Atlas project (44). Using these enhancer and gene expression data, we constructed tissue-specific GRNs using ANANSE, and then calculated the TF influence scores for each of the tissues when taking the combination of all other tissues as the source tissue (Supplementary Table S7). We clustered the 18 tissues based on the correlation between TF influence scores us-

ing hierarchical clustering, showing that the influence score captures regulatory similarities and differences between tissues (Figure 6A and Supplementary Figure S10). For example, the esophagus and the skin cluster together, as these tissues are composed mostly of stratified squamous epithelial cells, and skeletal muscle and heart tissue are clustered together as both tissues contain striated muscle tissues.

For all studied tissues, we have provided a rich resource of key TFs of each tissue, with a list of top ten key TFs (Figure 6B). Many TFs in this list are known to play important functions for specific tissues, e.g. ELF3 and KLF5 for stomach, colon, and small intestine (122,123); TFAP2A, TFAP2C, TP63 and GRHL2 for the skin and esophagus (15,124,125); SOX2, SOX8 and OLIG1/2 for brain (126–128); and SPI1 for lung, spleen and bone marrow (129) (Figure 6A).

In summary, using ANANSE, we predicted key TFs for 18 human normal tissues. Many of these predicted TFs correlate well with the known literature of these tissues. In addition, the predicted key TFs in each tissue also provide us a rich resource to unveil TFs with novel functions in specific tissues.

## DISCUSSION

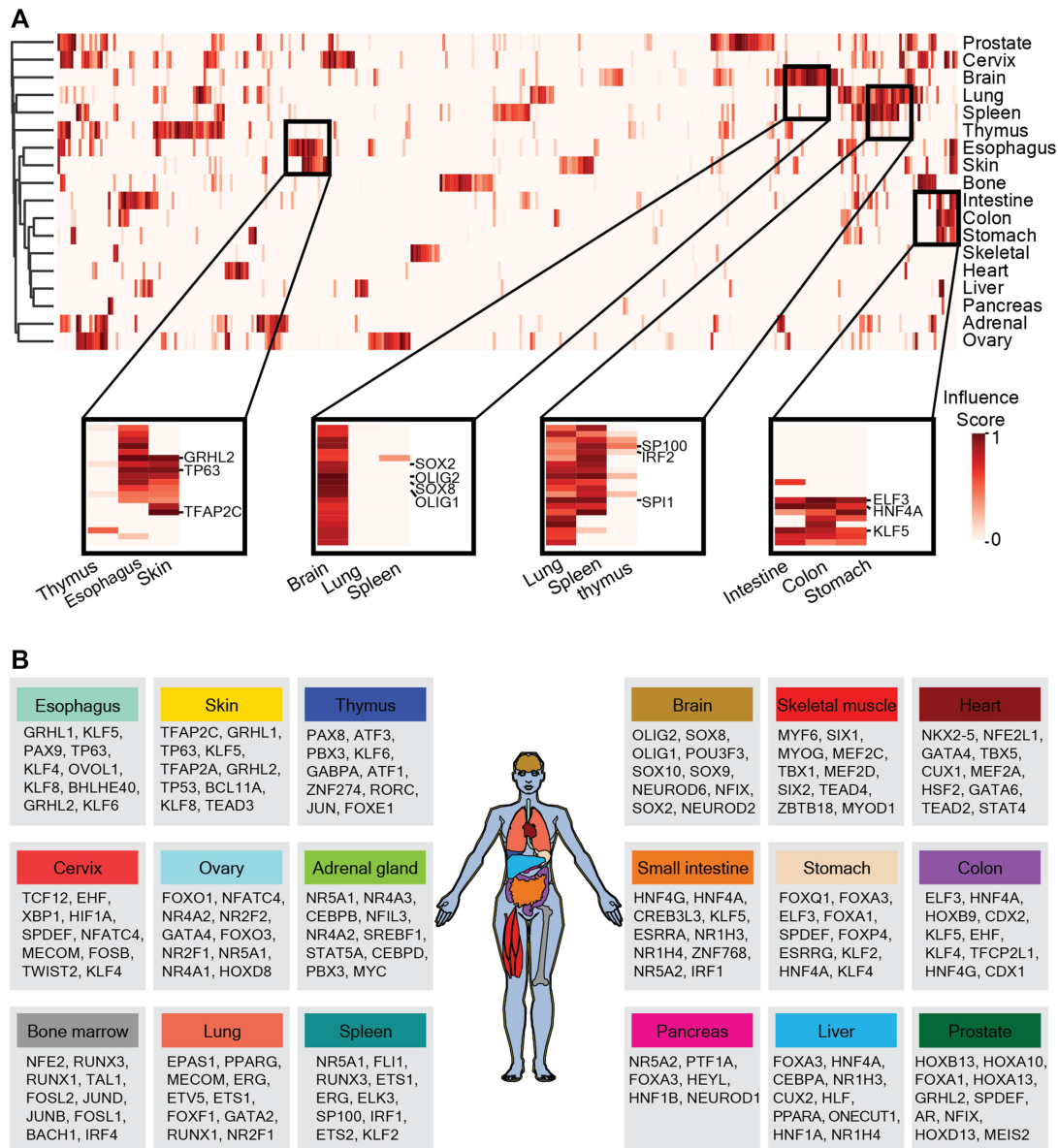
Lineage specification and cell fate determination are critical processes during development. They are necessary to form the diversity of cell types that are organized into organs and tissues. TFs form a central component in the regulatory networks that control lineage choice and differentiation. Indeed, cell fate can be switched *in vitro* through manipulation of TF expression (5,110–121). However, the regulatory factors that determine cell identity remain unknown for many cell types. To address this issue, we developed ANANSE, a new computational method to predict the key TFs that regulate cellular fate changes.

We establish TF binding networks for each cell type by leveraging genome-wide, cell type-specific enhancer signals from ATAC-seq, H3K27ac ChIP-seq and TF motif data. ANANSE takes a two-step approach. First, TF binding is imputed for all enhancers using a simple supervised logistic classifier. In contrast to existing methods that aim to predict binding by training TF-specific models (130–132), we also used a more general model that can be applied to all TFs. Logically, our simple model will not be as accurate as complex models, trained for specific TFs. Indeed, the PR AUC and ROC AUC of the ANANSE binding model are lower for a factor such as CTCF than the current state-of-the-art in supervised prediction, as illustrated by the comparison with Virtual ChIP-seq (104). However, the advantage of our model is that it can predict binding for every TF as long as its motif is known. In addition, it can be used for factors for which there is no training data available, and it can also be applied to non-model organisms that lack comprehensive ChIP-seq assays. Our benchmarks show that it performs significantly better than using the motif score alone or enhancer activity alone and that it outperformed a more stringent average ChIP-seq baseline.

Second, we summarized the imputed TF signals per gene, using a distance-weighted decay function (84), and combined this measure with TF activity and TF and target

gene expression to infer cell type-specific GRNs. In general, there is a lack of gold standards to evaluate cell type-specific GRNs. We used multiple orthogonal types of benchmarks: databases of known, experimentally identified TF–gene interactions, gene expression after TF perturbations, and functional enrichment using Gene Ontology annotation. The databases with known interactions that we used contain only a fraction of true regulatory interactions, and therefore this benchmark is affected by a large fraction of false negatives. All our benchmark evaluations demonstrate that ANANSE significantly enriches for true regulatory interactions. However, it also highlights that GRN inference is far from a solved problem. The PR AUC values are low, as is generally the case in eukaryotic GRN inference (133). Our comparison with several established gene expression-based GRN approaches, such as PANDA (29), GRNBoost2 (similar to GENIE3) (27,108) and ARACNE (26) shows that these methods also result in low PR AUC values. In addition to the improved performance, ANANSE has another clear benefit. While most GRN inference methods need a large collection of samples, ANANSE uses only two or three genome-wide measurements as input: gene expression and enhancer activity (H3K27ac ChIP-seq and/or ATAC-seq).

In contrast to previous approaches, our method takes advantage of TF binding in enhancers, instead of only gene expression differences or TF binding to proximal promoters. This resulted in significantly improved performance, as benchmarked for the GRN inference, as well as on experimentally validated *trans*-differentiation protocols. It has been previously shown that cell type-specific regulation is much better captured by enhancers as compared to promoter-proximal regulatory elements. For instance, TF binding and chromatin accessibility in distal elements better reflect the cell type identity of hematopoietic lineages than in promoters (134,135). Many important transcriptional regulators mainly bind at regulatory regions that are not proximal to the promoter. Indeed, our analysis of the genomic binding distribution of ~300 human TFs showed that cell type-specific TFs bind in enhancer regions more often than TFs that are more widely expressed (Figure 1C). Therefore, we reasoned that TF binding at enhancers would be essential to model cell fate and lineage decisions. We tested the application of the networks inferred by ANANSE to human *in vitro trans*-differentiation approaches. Earlier work showed that computational algorithms allow characterization of cellular fate transitions and rational prioritization of TF candidates for *trans*-differentiation (21,23,25). We implemented a network-based approach to prioritize TFs that determine cell fate changes. Using a collection of known, experimentally validated *trans*-differentiation protocols, we demonstrated that ANANSE consistently outperforms other published approaches. This means that cellular trajectories can be characterized using ANANSE to identify the TFs that are involved in cell fate changes. In comparison with a promoter-based approach, we show that using enhancer-based regulatory information contributes significantly to this increased performance (Figure 5A). One noticeable example is the *trans*-differentiation from fibroblasts or mesenchymal cells to keratinocytes. In current experimentally validated *trans*-differentiation methods, the



**Figure 6.** Applying ANANSE to expression data of human tissues to identify key transcription factors. (A) Heatmap of the predicted influence scores of all TFs using ANANSE on data from 18 human tissues. The color in the heatmap indicates the relative influence score, from low to high. The four small heatmaps highlighted below show important TFs in related tissues. (B) The top 10 key TFs of 18 tissues inferred by ANANSE. The color of the tissue is consistent with the tissue name in the box. The order of TF of each tissue is based on the influence score of the TF ranked from high to low.

epithelial master regulator TP63 is essential for establishing the keratinocyte cell fate (117,136). However, TP63 was not predicted in most of the previously published computational methods (21,23,25). One plausible explanation is that TP63 is a TF for specific epithelial cells and tissues and it binds predominantly (87%) to enhancers (15,31–33), whereas previous computational tools do not take enhancer properties into consideration.

We used ANANSE to identify tissue-specific TFs for different human tissues. We predicted the top 10 key TFs for all studied tissues. Many TFs in this list are known for important functions in these specific tissues. For example, some NK homeodomain, GATA, and T-box TFs are found in normal cardiac development, which have impor-

tant functions during heart specification, patterning, and differentiation (137–139). Many TFs of the SOX family are known to be critical for neural system development in brain tissue (126,127). The gastrointestinal tract tissues share a number of high influence score TFs such as ELF3, KLF5 and HNF4A, which play roles in stomach, colon, and small intestine development, and are consistent with the current research on gastrointestinal tract tissues (Figure 6A) (57,140,141). ELF3 is important in intestinal morphogenesis, homeostasis, and disease (57). The *Klf5* deletion in mouse leads to intestine epithelial damage and a reduction of colon proliferative crypt cells (141). Our analysis showed that TP63, TFAP2A, TFAP2C and GRHL1 are common important TFs in the skin and esophagus (Fig-

ure 6B). The function of these TFs has been well studied in the skin. TP63 is one of the TFs that is important in both skin and esophagus development (15,117,142). TP63 and TFAP2A have been used in *in vivo* reprogramming of wound-resident cells to generate skin epithelial tissue (117). Both TFAP2A and TFAP2C are required for proper early morphogenesis and development as well as terminal differentiation of the skin epidermis (143–145). GRHL1 is important for the functioning of the epidermis. Grhl1 knockout mice exhibit palmoplantar keratoderma, impaired hair anchoring, and desmosomal abnormalities (125). It would be interesting to investigate what roles they play in esophagus. PAX9 regulates squamous cell differentiation and carcinogenesis in the oro-oesophageal epithelium (146). Although not all predicted TFs are known to have an important role in specific tissues, further research is warranted. The TFs in the TF atlas predicted by ANANSE may also be good candidates for studying tissue development and engineering in regenerative medicine.

Another large benefit of the model that we implemented in ANANSE is the wide applicability. The source code of ANANSE is publicly available under a liberal license. ANANSE does not depend on large collections of reference data and it is straightforward to run on new data, such as different cell types or even species. The types of data required for this analysis are the following: gene expression data (RNA-seq) and genome-wide assays of enhancer activity. The enhancer data can be ATAC-seq, H3K27ac ChIP-seq or a combination of them, which is relatively easily obtained, not only in human cell types or in common model species, but also often in non-model species (147). The predictions of ANANSE, represented by TF binding, gene regulatory networks and TFs ranked by their influence score, are useful to study gene regulatory principles in a wide variety of contexts. While we used *trans*-differentiation experiments to benchmark the TF influence score, the utility of ANANSE is not limited to these types of experiments. It can also be used to study differentiation, cell type-specific gene regulation and developmental processes.

We also acknowledge limitations in our approach. In ANANSE, we link enhancer regions to genes on the basis of distance. For each TF and gene interaction pair, ANANSE only considers TF binding information located at most 100 kb up and downstream of the corresponding gene. Although data from a recent CRISPR enhancer interference screen showed that genomic distance is largely informative in predicting enhancer-target interactions (148), this approach may be limited when applying to genes regulated through ultra-long range regulation or through less abundant inter-chromosomal contacts (149). This limitation of our method can potentially be addressed using chromosome conformation capture techniques (3C) (150) or other adaptations as circular 3C (4C) (151,152), chromosome conformation capture carbon copy (5C) (153), chromatin immunoprecipitation using PET (ChIA-PET) (154) and Hi-C (155). However, these types of data are currently only available for a limited number of cell types, therefore incorporation of topology data would limit the broad utility and application of our approach. Another limitation is that the current implementation of ANANSE focuses on activating transcription factors. During cell differentiation and repro-

gramming, other factors such as transcriptional repressors and chromatin modifying enzymes also play an important role. These are currently not considered in ANANSE. Finally, similar to other genome-wide gene regulatory network inference methods, the performance of ANANSE is not yet optimal. While our benchmarks (Figure 4 A and B, Supplementary Figure S2) indicate that ANANSE improves upon other methods, it is clear that there is still much progress to be made in genome-wide GRN inference.

## CONCLUSION

Here, we presented ANANSE, a computational tool for (i) transcription factor binding prediction, (ii) gene regulatory network inference and (iii) efficient prediction of TFs in cell fate determination. It outperforms other published methods in GRN inference and in predicting TFs that can induce *trans*-differentiation. In addition, it is open source, freely available and can be easily applied to other cell types and in any species. In summary, ANANSE exploits the powerful impact enhancers have on gene regulatory networks, and it provides insights into TF mediated regulatory mechanisms underlying cell fate determination and development.

## DATA AVAILABILITY

ANANSE source code is available from <https://github.com/vanheeringen-lab/ANANSE>. Jupyter notebooks for supporting analyses are provided at <https://github.com/vanheeringen-lab/ANANSE-manuscript>. GRNs inferred with ANANSE are available from Zenodo (tissue-specific GRNs <https://doi.org/10.5281/zenodo.4814016>; cell type-specific GRNs <https://doi.org/10.5281/zenodo.4809062>). Tissue-specific GRNs inferred with GRNBoost2 are available from Zenodo (<https://doi.org/10.5281/zenodo.4814015>).

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## ACKNOWLEDGEMENTS

This study used data provided by the ENCODE consortium (<https://www.encodeproject.org/>) and the NIH Roadmap Epigenomics Consortium (<http://nihroadmap.nih.gov/epigenomics/>). We would like to thank Jos Smits, Branco Heuts, Jori de Leuw and Seline van den Oever for providing critical input during the development of ANANSE.

## FUNDING

Chinese Scholarship Council [201606230213 to Q.X.]; Netherlands Organization for Scientific Research [NWO grant 016.Vidi.189.081 to S.J.v.H.]; US National Institutes of Health [NICHD, R01HD069344 to G.J.C.V.]. Funding for open access charge: Netherlands Organization for Scientific Research.

*Conflict of interest statement.* None declared.



## REFERENCES

1. Davis, R.L., Weintraub, H. and Lassar, A.B. (1987) Expression of a single transfected cDNA converts fibroblasts to myoblasts. *Cell*, **51**, 987–1000.
2. Jopling, C., Boue, S. and Izpisua Belmonte, J.C. (2011) Dedifferentiation, transdifferentiation and reprogramming: three routes to regeneration. *Nat. Rev. Mol. Cell Biol.*, **12**, 79–89.
3. Pang, Z.P., Yang, N., Vierbuchen, T., Ostermeier, A., Fuentes, D.R., Yang, T.Q., Citri, A., Sebastiano, V., Marro, S., Südhof, T.C. *et al.* (2011) Induction of human neuronal cells by defined transcription factors. *Nature*, **476**, 220–223.
4. Stadhouders, R., Filion, G.J. and Graf, T. (2019) Transcription factors and 3D genome conformation in cell-fate decisions. *Nature*, **569**, 345–354.
5. Takahashi, K., Tanabe, K., Ohnuki, M., Narita, M., Ichisaka, T., Tomoda, K. and Yamanaka, S. (2007) Induction of pluripotent stem cells from adult human fibroblasts by defined factors. *Cell*, **131**, 861–872.
6. Vierbuchen, T., Ostermeier, A., Pang, Z.P., Kokubu, Y., Südhof, T.C. and Wernig, M. (2010) Direct conversion of fibroblasts to functional neurons by defined factors. *Nature*, **463**, 1035–1041.
7. Lambert, S.A., Jolma, A., Campitelli, L.F., Das, P.K., Yin, Y., Albu, M., Chen, X., Taipale, J., Hughes, T.R. and Weirauch, M.T. (2018) The Human Transcription Factors. *Cell*, **172**, 650–665.
8. Vaquerizas, J.M., Kummerfeld, S.K., Teichmann, S.A. and Luscombe, N.M. (2009) A census of human transcription factors: function, expression and evolution. *Nat. Rev. Genet.*, **10**, 252–263.
9. Davidson, E.H. (2010) Emerging properties of animal gene regulatory networks. *Nature*, **468**, 911–920.
10. Tegner, J. and Björkegren, J. (2007) Perturbations to uncover gene networks. *Trends Genet.*, **23**, 34–41.
11. Wilkinson, A.C., Nakauchi, H. and Göttgens, B. (2017) Mammalian transcription factor networks: recent advances in interrogating biological complexity. *Cell Syst.*, **5**, 319–331.
12. Iwafuchi-Doi, M. and Zaret, K.S. (2016) Cell fate control by pioneer transcription factors. *Development*, **143**, 1833–1837.
13. Peñalosa-Ruiz, G., Bright, A.R., Mulder, K.W. and Veenstra, G.J.C. (2019) The interplay of chromatin and transcription factors during cell fate transitions in development and reprogramming. *Biochim. Biophys. Acta Gene Regul. Mech.*, **1862**, 194407.
14. Buschbeck, M. and Hake, S.B. (2017) Variants of core histones and their roles in cell fate decisions, development and cancer. *Nat. Rev. Mol. Cell Biol.*, **18**, 299–314.
15. Qu, J., Tanis, S.E., Smits, J.P., Kouwenhoven, E.N., Oti, M., van den Bogaard, E.H., Logie, C., Stunnenberg, H.G., van Bokhoven, H. and Mulder, K.W. (2018) Mutant p63 affects epidermal cell identity through rewiring the enhancer landscape. *Cell Rep.*, **25**, 3490–3503.
16. Reik, W., Dean, W. and Walter, J. (2001) Epigenetic reprogramming in mammalian development. *Science (New York, N.Y.)*, **293**, 1089–1093.
17. Heinaniemi, M., Nykter, M., Kramer, R., Wienecke-Baldacchino, A., Sinkkonen, L., Zhou, J.X., Kreisberg, R., Kauffman, S.A., Huang, S. and Shmulevich, I. (2013) Gene-pair expression signatures reveal lineage control. *Nat. Methods*, **10**, 577–583.
18. Lang, A.H., Li, H., Collins, J.J. and Mehta, P. (2014) Epigenetic landscapes explain partially reprogrammed cells and identify key reprogramming genes. *PLoS Comput. Biol.*, **10**, e1003734.
19. Roost, M.S., van Iperen, L., Ariyurek, Y., Buermans, H.P., Arindrarto, W., Devalla, H.D., Passier, R., Mummery, C.L., Carloti, F., de Koning, E.J. *et al.* (2015) KeyGenes, a tool to probe tissue differentiation using a human fetal transcriptional atlas. *Stem Cell Rep.*, **4**, 1112–1124.
20. D'Alessio, A.C., Fan, Z.P., Wert, K.J., Baranov, P., Cohen, M.A., Saini, J.S., Cohick, E., Charniga, C., Dadon, D. and Hannett, N.M. (2015) A systematic approach to identify candidate transcription factors that control cell identity. *Stem Cell Rep.*, **5**, 763–775.
21. Cahan, P., Li, H., Morris, S.A., Da Rocha, E.L., Daley, G.Q. and Collins, J.J. (2014) CellNet: network biology applied to stem cell engineering. *Cell*, **158**, 903–915.
22. Hartmann, A., Okawa, S., Zaffaroni, G. and Del Sol, A. (2018) SeesawPred: a web application for predicting cell-fate determinants in cell differentiation. *Sci. Rep.*, **8**, 13355.
23. Morris, S.A., Cahan, P., Li, H., Zhao, A.M., San Roman, A.K., Shivdasani, R.A., Collins, J.J. and Daley, G.Q. (2014) Dissecting engineered cell types and enhancing cell fate conversion via CellNet. *Cell*, **158**, 889–902.
24. Alvarez, M.J., Shen, Y., Giorgi, F.M., Lachmann, A., Ding, B.B., Ye, B.H. and Califano, A. (2016) Functional characterization of somatic mutations in cancer using network-based inference of protein activity. *Nat. Genet.*, **48**, 838–847.
25. Rackham, O.J., Firas, J., Fang, H., Oates, M.E., Holmes, M.L., Knaupp, A.S., Suzuki, H., Nefzger, C.M., Daub, C.O. and Shin, J.W. (2016) A predictive computational framework for direct reprogramming between human cell types. *Nat. Genet.*, **48**, 331.
26. Margolin, A.A., Nemenman, I., Basso, K., Wiggins, C., Stolovitzky, G., Dalla Favera, R. and Califano, A. (2006) ARACNE: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinformatics*, **7**, S7.
27. Huynh-Thu, V.A., Irrthum, A., Wehenkel, L. and Geurts, P. (2010) Inferring regulatory networks from expression data using tree-based methods. *PLoS One*, **5**, e12776.
28. Marbach, D., Roy, S., Ay, F., Meyer, P.E., Candeias, R., Kahveci, T., Bristow, C.A. and Kellis, M. (2012) Predictive regulatory models in *Drosophila melanogaster* by integrative inference of transcriptional networks. *Genome Res.*, **22**, 1334–1349.
29. Glass, K., Huttenhower, C., Quackenbush, J. and Yuan, G.-C. (2013) Passing messages between biological networks to refine predicted interactions. *PLoS One*, **8**, e64832.
30. Janky, R., Verfaillie, A., Imrichová, H., Van de Sande, B., Standaert, L., Christiaens, V., Hulselmans, G., Hertens, K., Naval Sanchez, M., Potier, D. *et al.* (2014) iRegulon: from a gene list to a gene regulatory network using large motif and track collections. *PLoS Comput. Biol.*, **10**, e1003731.
31. Andersson, R., Gebhard, C., Miguel-Escalada, I., Hoof, I., Bornholdt, J., Boyd, M., Chen, Y., Zhao, X., Schmid, C., Suzuki, T. *et al.* (2014) An atlas of active enhancers across human cell types and tissues. *Nature*, **507**, 455–461.
32. Bulger, M. and Groudine, M. (2011) Functional and mechanistic diversity of distal transcription enhancers. *Cell*, **144**, 327–339.
33. Spitz, F. and Furlong, E.E.M. (2012) Transcription factors: from enhancer binding to developmental control. *Nat. Rev. Genet.*, **13**, 613–626.
34. Davis, C.A., Hitz, B.C., Sloan, C.A., Chan, E.T., Davidson, J.M., Gabdank, I., Hilton, J.A., Jain, K., Baymuradov, U.K., Narayanan, A.K. *et al.* (2018) The Encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic Acids Res.*, **46**, D794–D801.
35. Valouev, A., Johnson, D.S., Sundquist, A., Medina, C., Anton, E., Batzoglou, S., Myers, R.M. and Sidow, A. (2008) Genome-wide analysis of transcription factor binding sites based on ChIP-Seq data. *Nat. Methods*, **5**, 829.
36. Luo, Z., Gao, X., Lin, C., Smith, E.R., Marshall, S.A., Swanson, S.K., Florens, L., Washburn, M.P. and Shilatifard, A. (2015) Zic2 is an enhancer-binding factor required for embryonic stem cell specification. *Mol. Cell*, **57**, 685–694.
37. Whyte, W.A., Orlando, D.A., Hnisz, D., Abraham, B.J., Lin, C.Y., Kagey, M.H., Rahl, P.B., Lee, T.I. and Young, R.A. (2013) Master transcription factors and mediator establish super-enhancers at key cell identity genes. *Cell*, **153**, 307–319.
38. Huang, M., Chen, Y., Yang, M., Guo, A., Xu, Y., Xu, L. and Koeffer, H.P. (2017) dbCoRC: a database of core transcriptional regulatory circuitries modeled by H3K27ac ChIP-seq signals. *Nucleic Acids Res.*, **46**, D71–D77.
39. Qin, Q., Fan, J., Zheng, R., Wan, C., Mei, S., Wu, Q., Sun, H., Brown, M., Zhang, J., Meyer, C.A. *et al.* (2020) Lisa: inferring transcriptional regulators through integrative modeling of public chromatin accessibility and ChIP-seq data. *Genome Biol.*, **21**, 32.
40. Wang, Z., Civelek, M., Miller, C.L., Sheffield, N.C., Guertin, M.J. and Zang, C. (2018) BART: a transcription factor prediction tool with query gene sets or epigenomic profiles. *Bioinformatics*, **34**, 2867–2869.
41. Chèneby, J., Gheorghe, M., Artufel, M., Mathelier, A. and Ballester, B. (2017) ReMap 2018: an updated atlas of regulatory regions from an integrative analysis of DNA-binding ChIP-seq experiments. *Nucleic Acids Res.*, **46**, D267–D275.

42. Yu, G., Wang, L.-G. and He, Q.-Y. (2015) ChIPseeker: an R/Bioconductor package for ChIP peak annotation, comparison and visualization. *Bioinformatics*, **31**, 2382–2383.
43. Korotkevich, G., Sukhov, V., Budin, N., Shpak, B., Artyomov, M.N. and Sergushichev, A. (2021) Fast gene set enrichment analysis. bioRxiv doi: <https://doi.org/10.1101/060012>, 01 February 2021, preprint: not peer reviewed.
44. Uhlén, M., Fagerberg, L., Hallström, B.M., Lindskog, C., Oksvold, P., Mardinoglu, A., Sivertsson, Å., Kampf, C., Sjöstedt, E. and Asplund, A. (2015) Tissue-based map of the human proteome. *Science*, **347**, 1260419.
45. Bernstein, B.E., Stamatoyannopoulos, J.A., Costello, J.F., Ren, B., Milosavljevic, A., Meissner, A., Kellis, M., Marra, M.A., Beaudet, A.L., Ecker, J.R. et al. (2010) The NIH roadmap epigenomics mapping consortium. *Nat. Biotechnol.*, **28**, 1045–1048.
46. Encode Project Consortium. (2012) An integrated encyclopedia of DNA elements in the human genome. *Nature*, **489**, 57–74.
47. Novakovic, B., Habibi, E., Wang, S.-Y., Arts, R.J., Davar, R., Megchelenbrink, W., Kim, B., Kuznetsova, T., Kox, M. and Zwaag, J. (2016)  $\beta$ -Glucan reverses the epigenetic state of LPS-induced immunological tolerance. *Cell*, **167**, 1354–1368.
48. Liu, Q., Jiang, C., Xu, J., Zhao, M.-T., Van Bortle, K., Cheng, X., Wang, G., Chang, H.Y., Wu, J.C. and Snyder, M.P. (2017) Genome-wide temporal profiling of transcriptome and open chromatin of early cardiomyocyte differentiation derived from hiPSCs and hESCs. *Circ. Res.*, **121**, 376–391.
49. Runge, J.S., Raab, J.R. and Magnuson, T. (2018) Identification of two distinct classes of the human INO80 complex genome-wide. *G3 (Bethesda)*, **8**, 1095–1102.
50. Buenostro, J.D., Corces, M.R., Lareau, C.A., Wu, B., Schep, A.N., Aryee, M.J., Majeti, R., Chang, H.Y. and Greenleaf, W.J. (2018) Integrated single-cell analysis maps the continuous regulatory landscape of human hematopoietic differentiation. *Cell*, **173**, 1535–1548.
51. Cho, S.W., Xu, J., Sun, R., Mumbach, M.R., Carter, A.C., Chen, Y.G., Yost, K.E., Kim, J., He, J., Nevins, S.A. et al. (2018) Promoter of lncRNA gene PVT1 is a tumor-suppressor DNA boundary element. *Cell*, **173**, 1398–1412.
52. Morris, J.A., Kemp, J.P., Youtlen, S.E., Laurent, L., Logan, J.G., Chai, R.C., Vulpescu, N.A., Forgetta, V., Kleinman, A., Mohanty, S.T. et al. (2019) An atlas of genetic influences on osteoporosis in humans and mice. *Nat. Genet.*, **51**, 258–266.
53. Oomen, M.E., Hansen, A.S., Liu, Y., Darzacq, X. and Dekker, J. (2019) CTCF sites display cell cycle-dependent dynamics in factor binding and nucleosome positioning. *Genome Res.*, **29**, 236–249.
54. Li, L., Wang, Y., Torkelson, J.L., Shankar, G., Pattison, J.M., Zhen, H.H., Fang, F., Duren, Z., Xin, J., Gaddam, S. et al. (2019) TFAP2C- and p63-dependent networks sequentially rearrange chromatin landscapes to drive human epidermal lineage commitment. *Cell Stem Cell*, **24**, 271–284.
55. Tchieu, J., Calder, E.L., Guttikonda, S.R., Gutzwiller, E.M., Aromolaran, K.A., Steinbeck, J.A., Goldstein, P.A. and Studer, L. (2019) NFIA is a gliogenic switch enabling rapid derivation of functional human astrocytes from pluripotent stem cells. *Nat. Biotechnol.*, **37**, 267.
56. Song, M., Yang, X., Ren, X., Maliskova, L., Li, B., Jones, I.R., Wang, C., Jacob, F., Wu, K., Traglia, M. et al. (2019) Mapping cis-regulatory chromatin contacts in neural cells links neuropsychiatric disorder risk variants to target genes. *Nat. Genet.*, **51**, 1252–1262.
57. Johnston, A.D., Simões-Pires, C.A., Thompson, T.V., Suzuki, M. and Gready, J.M. (2019) Functional genetic variants can mediate their regulatory effects through alteration of transcription factor binding. *Nat. Commun.*, **10**, 3472.
58. Soares, E., Xu, Q., Li, Q., Qu, J., Zheng, Y., Raeven, H.H., Brandao, K.O., Petit, I., van den Akker, W.M. and van Heeringen, S.J. (2019) Single-cell RNA-seq identifies a reversible mesodermal activation in abnormally specified epithelia of p63 EEC syndrome. *Proc. Natl. Acad. Sci. U.S.A.*, **116**, 17361–17370.
59. Martone, J., Lisi, M., Castagnetti, F., Rosa, A., Di Carlo, V., Blanco, E., Setti, A., Mariani, D., Colantoni, A., Santini, T. et al. (2020) Trans-generational epigenetic regulation associated with the amelioration of Duchenne Muscular Dystrophy. *EMBO Mol. Med.*, **12**, e12063.
60. Grubert, F., Srivas, R., Spacek, D.V., Kasowski, M., Ruiz-Velasco, M., Sinnott-Armstrong, N., Greenside, P., Narasimha, A., Liu, Q., Geller, B. et al. (2020) Landscape of cohesin-mediated chromatin loops in the human genome. *Nature*, **583**, 737–743.
61. Segura-Bayona, S., Villamor-Paya, M., Attolini, C.S., Koenig, L.M., Sanchiz-Calvo, M., Boulton, S.J. and Stracker, T.H. (2020) Toused-like kinases suppress innate immune signaling triggered by alternative lengthening of telomeres. *Cell Rep.*, **32**, 107983.
62. Liu, Q., Zaba, L.C., Satpathy, A.T., Longmire, M., Zhang, W., Li, K., Granja, J., Guo, C., Lin, J., Li, R. et al. (2020) Chromatin accessibility landscapes of skin cells in systemic sclerosis nominate dendritic cells in disease pathogenesis. *Nat. Commun.*, **11**, 5843.
63. Granja, J.M., Corces, M.R., Pierce, S.E., Bagdatli, S.T., Choudhry, H., Chang, H.Y. and Greenleaf, W.J. (2021) ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis. *Nat. Genet.*, **53**, 403–411.
64. Sloan, C.A., Chan, E.T., Davidson, J.M., Malladi, V.S., Strattan, J.S., Hitz, B.C., Gabbank, I., Narayanan, A.K., Ho, M., Lee, B.T. et al. (2016) ENCODE data at the ENCODE portal. *Nucleic Acids Res.*, **44**, D726–D732.
65. van der Sande, M., Frölich, S., Smits, J. and van Heeringen, S.J. (2020) seq2science (Version v0.3.1). <https://doi.org/10.5281/zenodo.3921913>.
66. van Heeringen, S.J. (2017) genomepy: download genomes the easy way. *J. Open Source Softw.*, **2**, 320.
67. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M. and Gingeras, T.R. (2013) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, **29**, 15–21.
68. Picard2019toolkit (2019) Broad Institute, GitHub repository.
69. Zhang, Y., Liu, T., Meyer, C.A., Eeckhoutte, J., Johnson, D.S., Bernstein, B.E., Nusbaum, C., Myers, R.M., Brown, M. and Li, W. (2008) Model-based analysis of ChIP-Seq (MACS). *Genome Biol.*, **9**, R137.
70. Li, Q., Brown, J.B., Huang, H. and Bickel, P.J. (2011) Measuring reproducibility of high-throughput experiments. *Ann. Appl. Stat.*, **5**, 1752–1779.
71. Chen, S., Zhou, Y., Chen, Y. and Gu, J. (2018) fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*, **34**, i884–i890.
72. Li, H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv doi: <https://arxiv.org/abs/1303.3997>, 26 May 2013, preprint: not peer reviewed.
73. Amemiya, H.M., Kundaje, A. and Boyle, A.P. (2019) The ENCODE blacklist: identification of problematic regions of the genome. *Sci. Rep.*, **9**, 9354.
74. Patro, R., Duggal, G., Love, M.I., Irizarry, R.A. and Kingsford, C. (2017) Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods*, **14**, 417–419.
75. Cunningham, F., Achuthan, P., Akanni, W., Allen, J., Amode, M.R., Armean, I.M., Bennett, R., Bhari, J., Billis, K., Boddus, S. et al. (2018) Ensembl 2019. *Nucleic Acids Res.*, **47**, D745–D751.
76. Sonesson, C., Love, M.I. and Robinson, M.D. (2015) Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. *F1000Res*, **4**, 1521.
77. Love, M.I., Huber, W. and Anders, S. (2014) Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.*, **15**, 550.
78. Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, **26**, 841–842.
79. Bruse, N. and Heeringen, S.J.v. (2018) GimmeMotifs: an analysis framework for transcription factor motif analysis. bioRxiv doi: <https://doi.org/10.1101/474403>, 20 November 2018, preprint: not peer reviewed.
80. van Heeringen, S.J. and Veenstra, G.J.C. (2010) GimmeMotifs: a de novo motif prediction pipeline for ChIP-sequencing experiments. *Bioinformatics*, **27**, 270–271.
81. van der Sande, M. and van Heeringen, S.J. (2020) qnorm (Version v0.6.1). <https://doi.org/10.5281/zenodo.4114608>.
82. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R. and Dubourg, V. (2011) Scikit-learn: machine learning in Python. *J. Mach. Learn. Res.*, **12**, 2825–2830.

83. Weirauch, M.T., Yang, A., Albu, M., Cote, A.G., Montenegro-Montero, A., Drewe, P., Najafabadi, H.S., Lambert, S.A., Mann, I. and Cook, K. (2014) Determination and inference of eukaryotic transcription factor sequence specificity. *Cell*, **158**, 1431–1443.
84. Wang, S., Zang, C., Xiao, T., Fan, J., Mei, S., Qin, Q., Wu, Q., Li, X., Xu, K. and He, H.H. (2016) Modeling cis-regulation with a compendium of genome-wide histone H3K27ac profiles. *Genome Res.*, **26**, 1417–1429.
85. Dask Development Team (2016) Dask: Library for dynamic task scheduling.
86. Stovner, E.B. and Sætrum, P. (2020) PyRanges: efficient comparison of genomic intervals in Python. *Bioinformatics*, **36**, 918–919.
87. Aibar, S., González-Blas, C.B., Moerman, T., Huynh-Thu, V.A., Imrichova, H., Hulselmans, G., Rambow, F., Marine, J.C., Geurts, P., Aerts, J. *et al.* (2017) SCENIC: single-cell regulatory network inference and clustering. *Nat. Methods*, **14**, 1083–1086.
88. Fang, L., Li, Y., Ma, L., Xu, Q., Tan, F. and Chen, G. (2021) GRNdb: decoding the gene regulatory networks in diverse human and mouse conditions. *Nucleic Acids Res.*, **49**, D97–D103.
89. Mercatelli, D., Lopez-Garcia, G. and Giorgi, F.M. (2020) corto: a lightweight R package for gene network inference and master regulator analysis. *Bioinformatics*, **36**, 3916–3917.
90. Holland, C.H., Tanevski, J., Perales-Patón, J., Gleixner, J., Kumar, M.P., Mereu, E., Joughin, B.A., Stegle, O., Lauffenburger, D.A., Heyn, H. *et al.* (2020) Robustness and applicability of transcription factor and pathway analysis tools on single-cell RNA-seq data. *Genome Biol.*, **21**, 36.
91. Liu, Z.-P., Wu, C., Miao, H. and Wu, H. (2015) RegNetwork: an integrated database of transcriptional and post-transcriptional regulatory networks in human and mouse. *Database*, **2015**, bav095.
92. Han, H., Cho, J.-W., Lee, S., Yun, A., Kim, H., Bae, D., Yang, S., Kim, C.Y., Lee, M. and Kim, E. (2017) TRRUST v2: an expanded reference database of human and mouse transcriptional regulatory interactions. *Nucleic Acids Res.*, **46**, D380–D386.
93. Liberzon, A., Subramanian, A., Pinchback, R., Thorvaldsdóttir, H., Tamayo, P. and Mesirov, J.P. (2011) Molecular signatures database (MSigDB) 3.0. *Bioinformatics*, **27**, 1739–1740.
94. Obayashi, T., Kagaya, Y., Aoki, Y., Tadaka, S. and Kinoshita, K. (2019) COXPRESdb v7: a gene coexpression database for 11 animal species supported by 23 coexpression platforms for technical evaluation and evolutionary inference. *Nucleic Acids Res.*, **47**, D55–D62.
95. Kuleshov, M.V., Jones, M.R., Rouillard, A.D., Fernandez, N.F., Duan, Q., Wang, Z., Koplev, S., Jenkins, S.L., Jagodnik, K.M., Lachmann, A. *et al.* (2016) Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.*, **44**, W90–W97.
96. Chen, E.Y., Tan, C.M., Kou, Y., Duan, Q., Wang, Z., Meirelles, G.V., Clark, N.R. and Ma'ayan, A. (2013) Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinformatics*, **14**, 128.
97. Uhlen, M., Oksvold, P., Fagerberg, L., Lundberg, E., Jonasson, K., Forsberg, M., Zwaalen, M., Kampf, C., Wester, K. and Hober, S. (2010) Towards a knowledge-based human protein atlas. *Nat. Biotechnol.*, **28**, 1248.
98. Maag, J.L.V. (2018) gganatogram: An R package for modular visualisation of anatograms and tissues based on ggplot2. *F1000Res*, **7**, 1576–1576.
99. Hsu, F., Kent, W.J., Clawson, H., Kuhn, R.M., Diekhans, M. and Haussler, D. (2006) The UCSC known genes. *Bioinformatics*, **22**, 1036–1046.
100. Kim, Y.J., Lim, H., Li, Z., Oh, Y., Kovlyagina, I., Choi, I.Y., Dong, X. and Lee, G. (2014) Generation of multipotent induced neural crest by direct reprogramming of human postnatal fibroblasts with a single transcription factor. *Cell Stem Cell*, **15**, 497–506.
101. Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. and Greenleaf, W.J. (2013) Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat. Methods*, **10**, 1213–1218.
102. Schreiber, J., Singh, R., Bilmes, J. and Noble, W.S. (2020) A pitfall for machine learning methods aiming to predict across cell types. *Genome Biol.*, **21**, 282.
103. ENCODE-DREAM (2017) ENCODE-DREAM in vivo Transcription Factor Binding Site Prediction Challenge.
104. Karimzadeh, M. and Hoffman, M.M. (2018) Virtual ChIP-seq: predicting transcription factor binding by learning from the transcriptome. bioRxiv doi: <https://doi.org/10.1101/168419>, 12 March 2019, preprint: not peer reviewed.
105. Vierstra, J., Lazar, J., Sandstrom, R., Halow, J., Lee, K., Bates, D., Diegel, M., Dunn, D., Neri, F., Haugen, E. *et al.* (2020) Global reference mapping of human transcription factor footprints. *Nature*, **583**, 729–736.
106. Balwierz, P.J., Pachkov, M., Arnold, P., Gruber, A.J., Zavolan, M. and van Nimwegen, E. (2014) ISMARA: automated modeling of genomic signals as a democracy of regulatory motifs. *Genome Res.*, **24**, 869–884.
107. Suzuki, H., Forrest, A.R., van Nimwegen, E., Daub, C.O., Balwierz, P.J., Irvine, K.M., Lassmann, T., Ravasi, T., Hasegawa, Y., de Hoon, M.J. *et al.* (2009) The transcriptional network that controls growth arrest and differentiation in a human myeloid leukemia cell line. *Nat. Genet.*, **41**, 553–562.
108. Moerman, T., Aibar Santos, S., Bravo González-Blas, C., Simm, J., Moreau, Y., Aerts, J. and Aerts, S. (2019) GRNBoost2 and Arboreto: efficient and scalable inference of gene regulatory networks. *Bioinformatics*, **35**, 2159–2161.
109. The Gene Ontology, C. (2019) The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic Acids Res.*, **47**, D330–D338.
110. Ciazzio, M., Giannelli, S., Valente, P., Lignani, G., Carissimo, A., Sessa, A., Colasante, G., Bartolomeo, R., Massimo, L. and Ferroni, S. (2015) Direct conversion of fibroblasts into functional astrocytes by defined transcription factors. *Stem Cell Rep.*, **4**, 25–36.
111. Fu, J.-D., Stone, N.R., Liu, L., Spencer, C.I., Qian, L., Hayashi, Y., Delgado-Olguin, P., Ding, S., Bruneau, B.G. and Srivastava, D. (2013) Direct reprogramming of human fibroblasts toward a cardiomyocyte-like state. *Stem Cell Rep.*, **1**, 235–247.
112. Ifkovits, J.L., Addis, R.C., Epstein, J.A. and Gearhart, J.D. (2014) Inhibition of TGFβ signaling increases direct conversion of fibroblasts to induced cardiomyocytes. *PLoS One*, **9**, e89678.
113. Nakamori, D., Akamine, H., Takayama, K., Sakurai, F. and Mizuguchi, H. (2017) Direct conversion of human fibroblasts into hepatocyte-like cells by ATF5, PROX1, FOXA2, FOXA3, and HNF4A transduction. *Sci. Rep.*, **7**, 16675.
114. Simeonov, K.P. and Uppal, H. (2014) Direct reprogramming of human fibroblasts to hepatocyte-like cells by synthetic modified mRNAs. *PLoS One*, **9**, e100134.
115. Huangfu, D., Maehr, R., Guo, W., Eijkelenboom, A., Snitow, M., Chen, A.E. and Melton, D.A. (2008) Induction of pluripotent stem cells by defined factors is greatly improved by small-molecule compounds. *Nat. Biotechnol.*, **26**, 795–797.
116. Yu, J., Vodyanik, M.A., Smuga-Otto, K., Antosiewicz-Bourget, J., Frane, J.L., Tian, S., Nie, J., Jonsdottir, G.A., Ruotti, V., Stewart, R. *et al.* (2007) Induced pluripotent stem cell lines derived from human somatic cells. *Science*, **318**, 1917–1920.
117. Kurita, M., Araoka, T., Hishida, T., O'Keefe, D.D., Takahashi, Y., Sakamoto, A., Sakurai, M., Suzuki, K., Wu, J. and Yamamoto, M. (2018) In vivo reprogramming of wound-resident cells generates skin epithelial tissue. *Nature*, **561**, 243.
118. Feng, R., Desbordes, S.C., Xie, H., Tillo, E.S., Pixley, F., Stanley, E.R. and Graf, T. (2008) PU.1 and C/EBPα/β convert fibroblasts into macrophage-like cells. *Proc. Natl. Acad. Sci. U.S.A.*, **105**, 6057–6062.
119. Xie, H., Ye, M., Feng, R. and Graf, T. (2004) Stepwise reprogramming of B cells into macrophages. *Cell*, **117**, 663–676.
120. Li, Y., Wang, Y., Yu, J., Ma, Z., Bai, Q., Wu, X., Bao, P., Li, L., Ma, D. and Liu, J. (2017) Direct conversion of human fibroblasts into osteoblasts and osteocytes with small molecules and a single factor, Runx2. bioRxiv doi: <http://dx.doi.org/10.1101/127480>, 14 April 2017, preprint: not peer reviewed.
121. Yamamoto, K., Kishida, T., Sato, Y., Nishioka, K., Ejima, A., Fujiwara, H., Kubo, T., Yamamoto, T., Kanamura, N. and Mazda, O. (2015) Direct conversion of human fibroblasts into functional osteoblasts by defined factors. *Proc. Natl. Acad. Sci. U.S.A.*, **112**, 6152–6157.
122. Jedlicka, P. and Gutierrez-Hartmann, A.J.H. (2008) Ets transcription factors in intestinal morphogenesis, homeostasis and disease. *Histol. Histopathol.*, **23**, 1417.

123. Katz, J.P., Perreault, N., Goldstein, B.G., Lee, C.S., Labosky, P.A., Yang, V.W. and Kaestner, K.H. (2002) The zinc-finger transcription factor Klf4 is required for terminal differentiation of goblet cells in the colon. *Development*, **129**, 2619–2628.
124. Dollé, P.J.N.r.s. (2009) Developmental expression of retinoic acid receptors (RARs). *Nuclear Receptor Signal.*, **7**, nrs.07006.
125. Wilanowski, T., Caddy, J., Ting, S.B., Hislop, N.R., Cerruti, L., Auden, A., Zhao, L.L., Asquith, S., Ellis, S., Sinclair, R. *et al.* (2008) Perturbed desmosomal cadherin expression in grainy head-like 1-null mice. *EMBO J.*, **27**, 886–897.
126. Bani-Yaghoob, M., Tremblay, R.G., Lei, J.X., Zhang, D., Zurakowski, B., Sandhu, J.K., Smith, B., Ribocco-Lutkiewicz, M., Kennedy, J., Walker, P.R. *et al.* (2006) Role of Sox2 in the development of the mouse neocortex. *Dev. Biol.*, **295**, 52–66.
127. Muto, A., Iida, A., Satoh, S. and Watanabe, S. (2009) The group E Sox genes Sox8 and Sox9 are regulated by Notch signaling and are required for Muller glial cell development in mouse retina. *Exp. Eye Res.*, **89**, 549–558.
128. Meijer, D.H., Kane, M.F., Mehta, S., Liu, H., Harrington, E., Taylor, C.M., Stiles, C.D. and Rowitch, D.H. (2012) Separated at birth? The functional and molecular divergence of OLIG1 and OLIG2. *Nat. Rev. Neurosci.*, **13**, 819–831.
129. Ohteki, T., Maki, C. and Koyasu, S. (2001) Overexpression of Bcl-2 differentially restores development of thymus-derived CD4<sup>+</sup>8<sup>+</sup>T cells and intestinal intraepithelial T cells in IFN-regulatory factor-1-deficient mice. *J. Immunol.*, **166**, 6509–6513.
130. Keilwagen, J., Posch, S. and Grau, J. (2019) Accurate prediction of cell type-specific transcription factor binding. *Genome Biol.*, **20**, 9.
131. Li, H., Quang, D. and Guan, Y. (2019) Anchor: trans-cell type prediction of transcription factor binding sites. *Genome Res.*, **29**, 281–292.
132. Quang, D. and Xie, X. (2019) FactorNet: A deep learning framework for predicting cell type specific transcription factor binding from nucleotide-resolution sequential data. *Methods*, **166**, 40–47.
133. Chen, S. and Mar, J.C. (2018) Evaluating methods of inferring gene regulatory networks highlights their lack of performance for single cell gene expression data. *BMC Bioinformatics*, **19**, 232–232.
134. Corces, M.R., Buenrostro, J.D., Wu, B., Greenside, P.G., Chan, S.M., Koenig, J.L., Snyder, M.P., Pritchard, J.K., Kundaje, A., Greenleaf, W.J. *et al.* (2016) Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nat. Genet.*, **48**, 1193–1203.
135. Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., Laslo, P., Cheng, J.X., Murre, C., Singh, H. and Glass, C.K. (2010) Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol. Cell*, **38**, 576–589.
136. Chen, Y., Mistry, D.S. and Sen, G.L. (2014) Highly rapid and efficient conversion of human fibroblasts to keratinocyte-like cells. *J. Invest. Dermatol.*, **134**, 335–344.
137. Bruneau, B.G. (2013) Signaling and transcriptional networks in heart development and regeneration. *Cold Spring Harb. Perspect. Biol.*, **5**, a008292.
138. Kathiriya, I.S., Nora, E.P. and Bruneau, B.G. (2015) Investigating the transcriptional control of cardiovascular development. *Circ. Res.*, **116**, 700–714.
139. Stefanovic, S. and Christoffels, V.M. (2015) GATA-dependent transcriptional and epigenetic control of cardiac lineage specification and differentiation. *Cell. Mol. Life Sci.*, **72**, 3871–3881.
140. Thompson, C.A., DeLaForest, A. and Battle, M.A. (2018) Patterning the gastrointestinal epithelium to confer regional-specific functions. *Dev. Biol.*, **435**, 97–108.
141. Nandan, M.O., Ghaleb, A.M., Liu, Y., Bialkowska, A.B., McConnell, B.B., Shroyer, K.R., Robine, S. and Yang, V.W. (2014) Inducible intestine-specific deletion of Krüppel-like factor 5 is characterized by a regenerative response in adult mouse colon. *Dev. Biol.*, **387**, 191–202.
142. Daniely, Y., Liao, G., Dixon, D., Linnoila, R.I., Lori, A., Randell, S.H., Oren, M. and Jetten, A.M. (2004) Critical role of p63 in the development of a normal esophageal and tracheobronchial epithelium. *Am J Physiol. Cell Physiol.*, **287**, C171–C181.
143. Budirahardja, Y., Tan, P.Y., Doan, T., Weisdepp, P. and Zaidel-Bar, R. (2016) The AP-2 transcription factor APTF-2 is required for neuroblast and epidermal morphogenesis in *Caenorhabditis elegans* embryogenesis. *PLoS Genet.*, **12**, e1006048.
144. Kousa, Y.A., Fuller, E. and Schutte, B.C. (2018) IRF6 and AP2A interaction regulates epidermal development. *J. Invest. Dermatol.*, **138**, 2578–2588.
145. Wang, X., Pasolli, H.A., Williams, T. and Fuchs, E. (2008) AP-2 factors act in concert with Notch to orchestrate terminal differentiation in skin epidermis. *J. Cell Biol.*, **183**, 37–48.
146. Xiong, Z., Ren, S., Chen, H., Liu, Y., Huang, C., Zhang, Y.L., Odera, J.O., Chen, T., Kist, R., Peters, H. *et al.* (2018) PAX9 regulates squamous cell differentiation and carcinogenesis in the oro-oesophageal epithelium. *J. Pathol.*, **244**, 164–175.
147. Villar, D., Berthelot, C., Aldridge, S., Rayner, T.F., Lukk, M., Pignatelli, M., Park, T.J., Deaville, R., Erichsen, J.T., Jasinska, A.J. *et al.* (2015) Enhancer evolution across 20 mammalian species. *Cell*, **160**, 554–566.
148. Fulco, C.P., Nasser, J., Jones, T.R., Munson, G., Bergman, D.T., Subramanian, V., Grossman, S.R., Anyoha, R., Doughty, B.R., Patwardhan, T.A. *et al.* (2019) Activity-by-Contact model of enhancer-promoter regulation from thousands of CRISPR perturbations. *Nat. Genet.*, **51**, 1664–1669.
149. Olivares-Chauvet, P., Mukamel, Z., Lifshitz, A., Schwartzman, O., Elkayam, N.O., Lubling, Y., Deikus, G., Sebra, R.P. and Tanay, A. (2016) Capturing pairwise and multi-way chromosomal conformations using chromosomal walks. *Nature*, **540**, 296–300.
150. Dekker, J., Rippe, K., Dekker, M. and Kleckner, N. (2002) Capturing chromosome conformation. *Science*, **295**, 1306–1311.
151. Simonis, M., Klous, P., Splinter, E., Moshkin, Y., Willemsen, R., de Wit, E., van Steensel, B. and de Laat, W. (2006) Nuclear organization of active and inactive chromatin domains uncovered by chromosome conformation capture-on-chip (4C). *Nat. Genet.*, **38**, 1348–1354.
152. Zhao, Z., Tavoosidana, G., Sjölander, M., Göndör, A., Mariano, P., Wang, S., Kanduri, C., Lezcano, M., Sandhu, K.S., Singh, U. *et al.* (2006) Circular chromosome conformation capture (4C) uncovers extensive networks of epigenetically regulated intra- and interchromosomal interactions. *Nat. Genet.*, **38**, 1341–1347.
153. Dostie, J., Richmond, T.A., Arnaout, R.A., Selzer, R.R., Lee, W.L., Honan, T.A., Rubio, E.D., Krumm, A., Lamb, J., Nusbaum, C. *et al.* (2006) Chromosome Conformation Capture Carbon Copy (5C): a massively parallel solution for mapping interactions between genomic elements. *Genome Res.*, **16**, 1299–1309.
154. Fullwood, M.J., Liu, M.H., Pan, Y.F., Liu, J., Xu, H., Mohamed, Y.B., Orlov, Y.L., Velkov, S., Ho, A., Mei, P.H. *et al.* (2009) An oestrogen-receptor- $\alpha$ -bound human chromatin interactome. *Nature*, **462**, 58–64.
155. Lieberman-Aiden, E., van Berkum, N.L., Williams, L., Imakaev, M., Ragoczy, T., Telling, A., Amit, I., Lajoie, B.R., Sabo, P.J., Dorschner, M.O. *et al.* (2009) Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science*, **326**, 289–293.