

Figure S2: Comparison of Tissue-Specific Elements. Related to Figure 2, Section S5, and Section S6. Percentage of (A) edges, (B) genes, and (C) TFs that were identified as specific in the tissue listed along the Y-axis, that are also identified as specific to the tissue listed along the X-axis. (D) Comparison of the tissue-specific edges identified using PANDA-networks to those that would have been identified using a network defined based on co-expression information. (E) Comparison of the tissue-specific edges identified using PANDA-networks to functional edges identified by GIANT-networks.

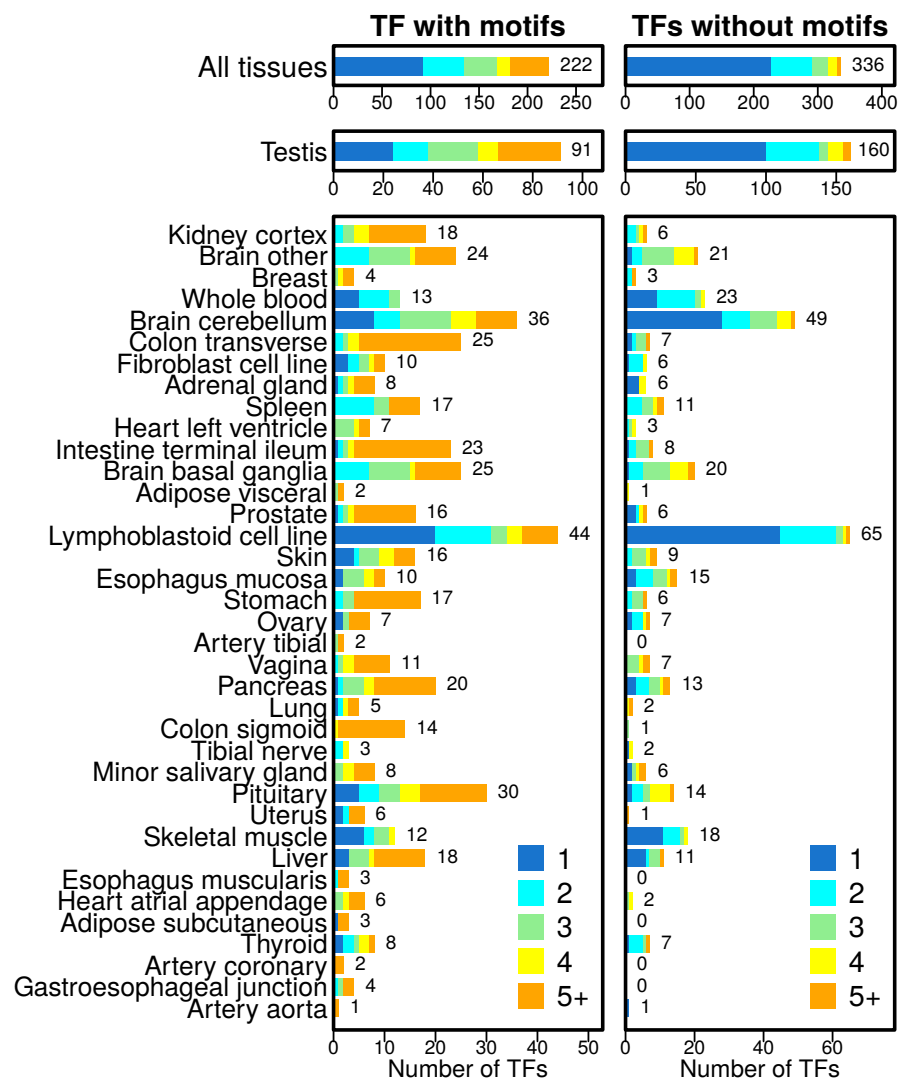


Figure S3: Evaluation of TFs without motifs. Related to Figure 2 and Section S4. Bar plots illustrating the number of TFs with (A) and without (B) DNA-binding motifs that were identified as “specific” to each of the 38 GTEx tissues. The total number of tissue-specific TFs identified for each tissue is shown to the right of each bar. Tissues are ordered on the number of edges specific to each tissue, as in Figure 2. Tissue-specificity for TFs was defined based on a TF having increased expression in one tissue compared to others, thus some TFs were identified as specific to multiple tissues. This multiplicity value is indicated by the color of the bars. We found that TFs with DNA-binding motifs have substantially higher multiplicity levels compared to TFs without motifs.

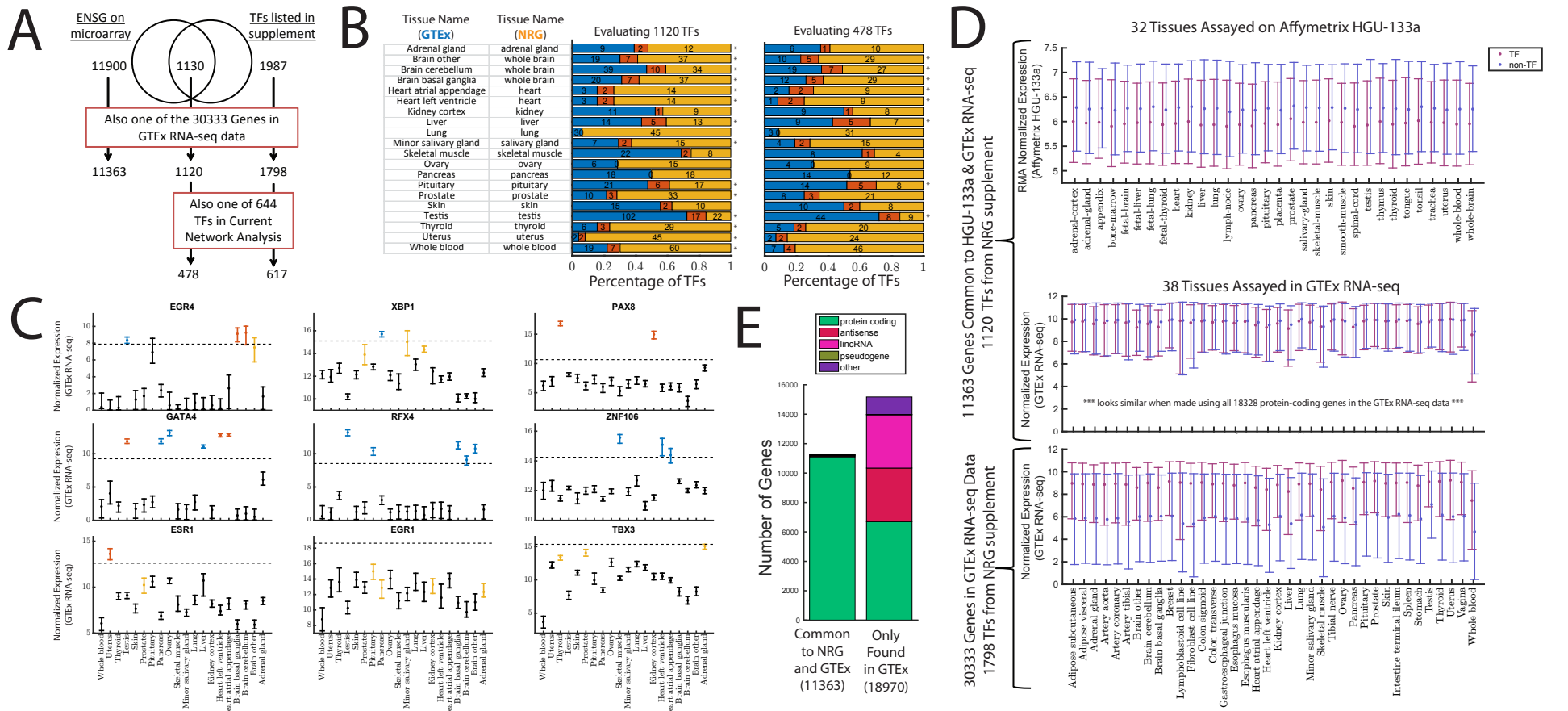


Figure S4: Analysis comparing the results of from a previous publication (NRG) with those obtained in this analysis using the GTEx RNA-seq data. Related to Figure 2 and Section S7. (A) An overview of the overlap in the genes included in the NRG gene expression data, the TFs included in the NRG supplemental data file, and how those sets overlap with the 30,333 genes in the normalized RNA-seq data we used in this analysis (see Section S1). (B) An analysis comparing the overlap of TFs identified as specific based on the NRG publication and those identified based on the GTEx data (see Section S4). An asterisk (*) indicates that the overlap is nominally significant ($p < 0.01$ by Fisher's exact test). (C) The distribution of expression values in the GTEx data for several example TFs. These TFs were chosen to illustrate a range of possibilities, including some overlap (*EGR4*, *GATA4*, *ESR1*), as well as opposing (*XBP1*), identical (*PAX8*), or distinct (*RFX4*, *ZNF106*, *EGR1*, *TBX3*) tissue-specific calls based on using either the NRG or the GTEx analysis. As there was little overlap between NRG and GTEx, the four plots with distinct tissue-specific calls are the most representative. (D) The expression of transcription factors versus non-transcription factor genes in both the NRG and GTEx expression data and using various criteria. (E) Information regarding the types of genes that are common between the set on the NRG microarray and in the GTEx RNA-seq data, and the types of genes that we have included in our GTEx expression analysis that were not on the NRG microarray.



Figure S5: Comparison of TF specificity based on expression versus network targeting. Related to Figure 5 and Section S9. Overlap between the TFs defined as tissue-specific based on their expression profile, versus those identified based on their differential targeting profile. All TFs that have tissue-specific differential targeting profiles can be found in Table S5.

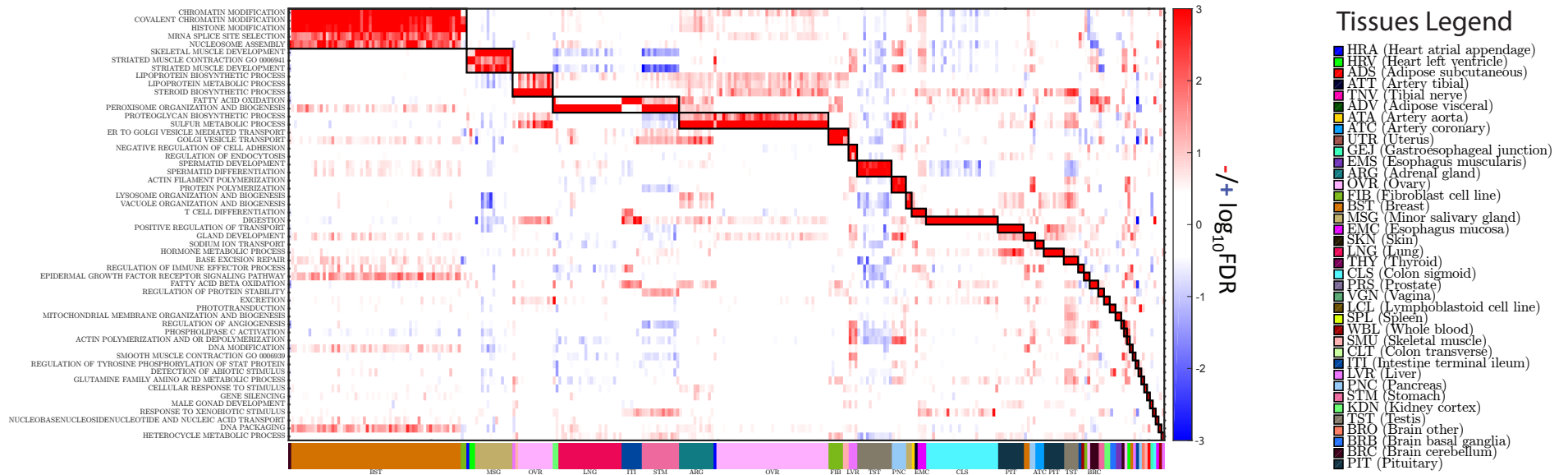


Figure S6: Illustration of the communities of GO terms and TF-tissue pairs that had three or fewer GO-term members. Related to Figure 5 and Section S9.

Distribution of Centrality Values (In Motif Prior Network)

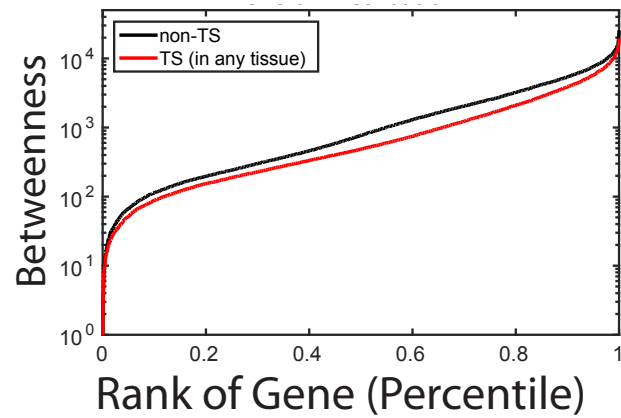
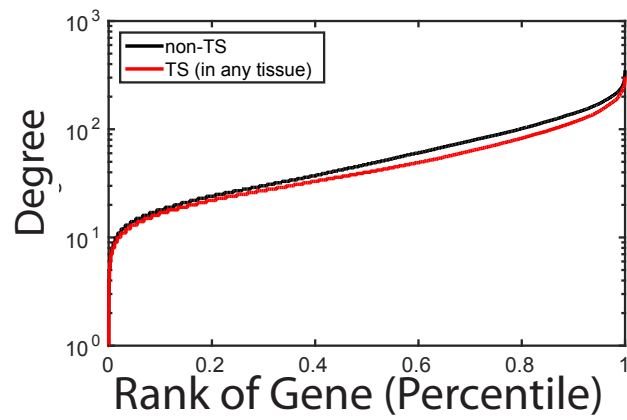


Figure S7: Centrality in the “prior” canonical regulatory network. Related to Figure 6 and Section S11. Distribution of the (A) in-degree and (B) betweenness centrality values of genes in the motif prior network used to seed the PANDA algorithm. Genes identified as tissue-specific are represented in the red line (all multiplicities considered), while those that are not identified as specific to any tissue are represented by the black line.

SUPPLEMENTAL EXPERIMENTAL PROCEDURES

S1. GTEx RNA-SEQ DATA

We downloaded the Genotype-Tissue Expression (GTEx) version 6.0 RNA-Seq data set (phs000424.v6.p1, 2015-10-05 released) from dbGaP (approved protocol #9112). GTEx release version 6.0 sampled 551 donors with phenotypic information and included 9,590 RNA-Seq assays. GTEx assayed expression in 30 tissue types, which were further divided into 53 tissue subregions (51 tissues and two derived cell lines) [1]. After removing tissues with very few samples (fewer than 15), we were left with 27 tissue types from 49 subregions. Using YARN (<http://bioconductor.org/packages/release/bioc/html/yarn.html>) we performed quality control, gene filtering, and normalization preprocessing. Briefly, we performed principal coordinate analysis (PCoA) using Y-chromosome genes to test for sample sex misidentification; we identified and removed GTEx-11ILO which was annotated as female but clustered with the males and was later confirmed to be an individual who underwent sex reassignment surgery (Kristin Ardlie, Broad Institute, private communication). We also used principal coordinate analysis on autosomal genes to group related body regions that had highly similar gene expression profiles. For example, skin samples from the lower leg (sun exposed) and from the suprapubic region (sun unexposed) shared gene expression profiles and were grouped as “skin,” while the transverse and descending colon were very different and were retained as distinct tissues. Gene expression data were then normalized using qsmooth [2] which performs a sparsity aware normalization that provides comparable expression profiles across all tissues and allowed us to retain genes that are expressed in only a single or small number of tissues. This preprocessing resulted in a dataset of 9,435 gene expression profiles assaying 30,333 genes in 38 tissues from 549 individuals. More detailed information on the normalization process and a complete description of the 38 final tissues and the associated samples are described elsewhere [3]. Consistent with GTEx, genes are denoted by their Ensembl IDs.

S2. REGULATORY NETWORK RECONSTRUCTION

We used the PANDA (Passing Attributes between Networks for Data Assimilation) network reconstruction algorithm [4] to estimate gene regulatory networks in each of the 38 GTEx tissues (see Section S1). PANDA incorporates regulatory information from three types of data: gene expression (used to create a co-expression network), protein-protein interaction, and a “prior” network based on mapping transcription factors to their putative target genes (used to initialize the algorithm).

Additional Gene Expression Data Processing: We filtered the normalized GTEx gene expression data (see above) to retain only the 30,243 genes that also had a significant motif-hit in their promoter region (see below). These genes were used when constructing our regulatory network models.

Prior Regulatory Network Based on Transcription Factor Motif Information: To create a “prior” regulatory network between transcription factors and genes, we downloaded Homo sapiens transcription factor motifs with direct/inferred evidence from the Catalog of Inferred Sequence Binding Preferences CIS-BP (<http://cisbp.cabr.utoronto.ca>, accessed: July 7, 2015). For each unique transcription factor, we selected the motif with the highest information content, resulting in a set of 695 motifs. We mapped these transcription factor position weight matrices (PWM) to the human genome (hg19) using FIMO [5] and retained highly significant matches ($p < 10^{-5}$) that occurred within the promoter regions of Ensembl genes (GRCh37.p13; annotations downloaded from <http://genome.ucsc.edu/cgi-bin/hgTables>, accessed: September 3, 2015); promoter regions were defined as $[-750, +250]$ around the transcription start site (TSS). After intersection to only include genes and transcription factors with expression data (see above) and at least one significant promoter hit, this process resulted in an initial map of potential regulatory interactions involving 644 transcription factors targeting 30,243 genes.

Prior Protein-Protein Interaction Network: We estimated an initial protein-protein interaction (PPI) network between all transcription factors (TFs) in our motif prior using interaction scores from StringDb v10 (<http://string-db.org>, accessed: October 27, 2015). PPI interaction scores were divided by 1,000 and self-interactions were set equal to one.

Reconstructing Networks using PANDA: For each of the 38 tissues, we used the GTEx gene expression data to calculate pairwise co-expression levels (based on Pearson correlation) between the 30,243 target genes. We then used PANDA to combine this information with the prior regulatory network and protein-protein interaction network. This produced 38 regulatory networks, one for each tissue, with edges predicted between 644 transcription factors and 30,243 target genes. PANDA returns complete, bipartite networks with edge weights similar to z-scores that represent the likelihood of a regulatory interaction. We transformed these z-scores to positive values using:

$$W_{ij}^{(t)} = \ln(e^{w_{ij}^{(t)}} + 1), \quad (\text{S1})$$

where $w_{ij}^{(t)}$ is the edge weight calculated by PANDA between a TF (i) and gene (j) in a particular tissue (t), and $W_{ij}^{(t)}$ is the transformed edge weight. These transformed edge weights are positive and so avoid issues related to calculating centrality measures on graph with negative edge weights (see Section S10).

S3. QUANTIFICATION OF TISSUE-SPECIFICITY VS GENERALITY OF NETWORK EDGES

Each of the 38 reconstructed PANDA networks contains scores, or “edge weights,” for every possible transcription factor-to-gene interaction (see Section S2). We used these edge weights to identify tissue-specific network edges. To do this, we compared the weight of an edge between a transcription factor (i) and a gene (j) in a particular tissue (t) to the median and interquartile range (IQR) of its weight across all 38 tissues:

$$s_{ij}^{(t)} = \frac{w_{ij}^{(t)} - \text{med}(w_{ij}^{(all)})}{\text{IQR}(w_{ij}^{(all)})}. \quad (\text{S2})$$

We then defined an edge with an edge specificity score $s_{ij}^{(t)} > N$ as specific to tissue t . We varied the cutoff N from 1 to 3, by steps of 0.25. Figure S1A shows the fraction of edges that are identified as tissue-specific at each cutoff. We selected a cutoff of $N = 2$ to define tissue-specific edges in order to be consistent with the cutoff used to define tissue-specific nodes (see Section S4). We also defined the “multiplicity” of an edge as:

$$m_{ij} = \sum_t (s_{ij}^{(t)} > N). \quad (\text{S3})$$

This value represents the number of tissues in which an edge is identified as specific. The overlap in edges identified as specific to each tissue can be found in Figure S2A.

S4. QUANTIFICATION OF TISSUE-SPECIFICITY VS GENERALITY OF NETWORK NODES

We wished to know if the tissue-specific edges were a direct reflection of the underlying gene expression data, or if the networks might be providing additional insight into the tissue-specific regulation of genes. Therefore, we identified tissue-specific network nodes (TFs and their target genes) by applying an analogous definition as we used to define tissue-specific edges to the GTEx gene expression data. We compared the median expression level $e_j^{(t)}$ of a gene (j) in a particular tissue (t), to the median and interquartile range of its expression across all samples:

$$s_j^{(t)} = \frac{\text{med}(e_j^{(t)}) - \text{med}(e_j^{(all)})}{\text{IQR}(e_j^{(all)})}. \quad (\text{S4})$$

We then defined a gene with gene specificity score $s_j^{(t)} > N$ as specific to tissue t . We varied the cutoff N from 1 to 3, by steps of 0.25. Figure S1B shows the fraction of tissue-specific genes identified at each cutoff. Based on this analysis, we selected a cutoff of $N = 2$ because with that cutoff approximately half of all genes are identified as tissue-specific. We also defined the “multiplicity” of a gene as:

$$m_j = \sum_t (s_j^{(t)} > N). \quad (\text{S5})$$

This value represents the number of tissues in which a gene is identified as specific. In Figure S1C we show some examples of non-tissue-specific and tissue-specific genes with different levels of multiplicity. We observe that the term “tissue-specific” is largely a misnomer. Many genes have a multiplicity greater than one, meaning that they are not actually “specific” to a particular tissue, but rather have a relatively higher level of expression in a subset of tissues compared to the others. Information regarding the tissue-specificity and multiplicity of genes can be found in Table S2. The overlap in genes identified as specific to each tissue can be found in Figure S2B.

To identify tissue-specific network nodes, other approaches, such as a differential expression analysis using Limma or ANOVA, could have been used. We decided to use the approach described in Equations S4–S5 for two reasons. First, this allowed us to link the tissue-specificity of network nodes to the specificity of network edges (Equations S2–S3). While multiple expression samples were available for each tissue (ranging from 36 to 779 samples, with a median of 210.5 samples), we only had one network available per tissue. We therefore could not use a statistical test that would compare groups of network edges. Second, approaches such as Limma or ANOVA assume that most genes are not differentially expressed between different conditions. In addition, these approaches assume normality. Because the GTEx expression data are not normally distributed, and because we observed global expression differences in expression between the different tissues (with large global expression differences in, for example, testis), these underlying null assumptions are likely not well founded.

Identifying Tissue-Specific Transcription Factors: Each of our network models includes information about the targeting profiles of 644 transcription factors (see Section S2). In analyzing tissue-specific transcription factors (Figure 2 in the main text) we focus on this subset of 644 transcription factors. However, other transcription factors that do not have a corresponding DNA-binding motif are included in our GTEx expression data. To evaluate these transcription factors we compared the genes in our GTEx expression data with a list of 1,987 genes that encode transcription factors as reported in a previous publication [6]. This list of 1,987 transcription factor genes includes 1,798 that have expression information in the normalized GTEx RNA-seq data (see Section S1) and 1,795

that are among the 30,243 target genes used in our network analysis (see Section S2). Of these 1,795 TFs, 1,178 do not have DNA-binding motifs, while 617 have a DNA-binding motif and are part of the set of 644 TFs we used in our regulatory prior. We compared the tissue-specificity and multiplicity levels of the 644 TFs with DNA-binding motifs and the 1,178 TFs without motifs (Figure S3). TFs with a motif were more likely to be identified as tissue-specific than TFs without motifs (34.5% versus 28.5%, Chi-squared test $p = 9.2 \cdot 10^{-3}$) and also tended to have higher multiplicity levels (Chi-squared test $p = 6.1 \cdot 10^{-14}$). Multiplicity levels of TF without motifs are comparable to those of all target genes (see Figure 2). Thus, while TFs with DNA-binding motifs are often shared among multiple tissues, TFs without such motifs are more often specific to only a single tissue, and appear to behave, at least in terms of their expression levels, more like non-TF genes. We note that although these differences may reflect biases in motif databases, they also potentially indicate that there are particular chemical and/or biological properties that have allowed certain transcription factors to be associated with known DNA-binding sequences.

S5. COMPARISON OF PANDA AND CORRELATION-BASED NETWORKS

Since co-expression networks have been widely used to analyze gene expression data, including in another network analysis of tissue-specificity in GTEx [7], we compared the tissue-specific edges defined based on PANDA-networks to those defined based on co-expression. For each of the 38 GTEx tissues analyzed we created co-expression networks by calculating the Pearson correlation between the 644 TFs and 30,243 genes included in our network model (see Section S2). We identified tissue-specific edges in these correlation-based networks using same protocol we used for genes and PANDA edges (Equation S2, with $N = 2$). When we compared the edges identified as tissue-specific using the correlation-based networks to those identified based on the PANDA-reconstructed regulatory networks we found very little overlap (Figure S2D).

This low level of overlap means that PANDA and Pearson correlation networks capture fundamentally different aspects of each tissue's gene expression program. The co-expression networks are based on measured expression correlations between TFs and their targets. In contrast, PANDA uses co-expression between target genes, not TFs and their targets. In particular, PANDA integrates target co-expression information with a prior regulatory network structure and TF-TF protein-protein interaction data, iteratively updating the likelihood of interactions between TFs and target genes based on shared patterns across these data.

We believe that PANDA more accurately captures tissue-specific regulatory processes. Indeed, when developing PANDA, we compared it to other methods, including co-expression networks, and found that the PANDA networks were better supported by confirmatory data, such as ChIP experiments [4]. Although no ChIP data are available for GTEx, PANDA does find biologically relevant associations that help elucidate the link between expression and tissue phenotype.

S6. COMPARISON OF PANDA AND FUNCTIONAL NETWORKS

We also systematically explored what, if any, information is shared between the PANDA tissue-specific regulatory networks we estimated from the GTEx data and other published tissue-specific functional networks [8]. In particular, we compared our networks to those estimated by GIANT (Genome-scale Integrated Analysis of gene Networks in Tissues; <http://giant.princeton.edu/>). This resource aims to capture tissue-specific functional interactions by using tissue-specific knowledge from the literature to selectively upweight particular datasets within a compendium of gene expression information and perform a tissue-ontology aware regularized Bayesian integration.

To begin, we identified which of the 144 distinct tissue-networks available on the GIANT web-resource are most likely to correspond to each of the 38 tissues we analyzed in the GTEx data. In some cases, we identified multiple tissue-networks available from GIANT to compare with a single tissue-network in our GTEx analysis. For example, in constructing the GTEx networks we combined gene expression samples from anatomically close regions of the brain when those samples were indistinguishable using principal component analysis (see Section S1 and [3]).

Next, we downloaded the top edges identified by GIANT, which represent interactions with a posterior probability greater than 0.1 in that tissue, for each of these identified corresponding tissues. As with the co-expression networks analyzed in Section S5, these GIANT networks are undirected, with edges extending between pairs of genes. Therefore, we next matched the nodes and dimensions between these functional networks and the GTEx regulatory networks. Specifically, of the 25,824 Entrez genes found across the GIANT networks, we identified 18,431 that were uniquely associated with the Ensembl and Gene Symbol annotations for the 30,243 genes that we had used to reconstruct the GTEx regulatory networks. This represents 91% of the unique Entrez gene annotations in the GTEx data (10,122 of the GTEx genes did not have a corresponding Entrez annotation). This set of 18,431 genes included 635 of our 644 transcription factors. We subsetted both the downloaded GIANT networks and our GTEx networks to only consider tissue-specific edges that extend from one of these 635 TFs to one of these 18,431 genes.

Finally, we determined the number of distinct and common edges between these subsetted GIANT and GTEx networks. We note that the number of top edges in the downloaded GIANT tissue-networks tended to be much greater than the number of tissue-specific edges we had identified from the GTEx PANDA networks. Even so, we found very little overlap between these functional networks and the gene regulatory networks (Figure S2E)

The low level of overlap between these different sets of tissue-specific networks is not surprising. PANDA and GIANT use very different approaches to build network models and, more importantly, the type of network each is designed to predict is fundamentally different. While PANDA aims to build a network representing the regulation of genes by transcription factors, GIANT's goal is to predict a functional network. Whereas edges in PANDA represent a directed interaction between a regulator and a target, edges in GIANT

networks represent the concordance across multiple independent lines of evidence that support an undirected association between two genes. Although both are “networks,” they represent very different types of biological relationships. This analysis highlights that the regulatory information contained in our GTEx regulatory networks is distinct from the type of information that is included in other tissue-network resources.

57. COMPARISON WITH A PREVIOUSLY PUBLISHED TISSUE-SPECIFIC TF RESOURCE

We also compared the transcription factors we identified as tissue-specific based on the GTEx expression data (see Section S4) with those reported as tissue-specific in a previous publication [6] (hereafter referred to as NRG, standing for the journal in which it was published: Nature Reviews Genetics) and which were used in other GTEx network evaluations [7]. The results of this analysis are shown in Figure S4.

To begin, we downloaded the gene expression data used for the calling of tissue-specific transcription factors in the NRG publication from the Gene Expression Omnibus (GSE1133). We RMA-normalized these expression data using the `justRMA()` function in the `affy` Version 1.52.0 library from Bioconductor in R and used a custom-CDF for the Affymetrix GeneChip HG-U133A array (`hgu133ahsengcdf_20.0.0`) [9] in order to normalize with respect to current Ensembl genes IDs. This RMA-normalized version of the expression data contained expression information for 11,900 different Ensembl genes across 158 total samples, 64 of which correspond to the “32 healthy major tissues and organs” used in the NRG analysis. 11,363 of the genes in this RMA-normalized NRG expression data set also appeared in the normalized GTEx data (see Section S1 and Figure S4A).

We next downloaded the supplemental data that accompanied the NRG manuscript. The “supplemental information S3” file contained information for 1,987 genes that encode transcription factors, including their “Ensembl gene IDs (release 51), HGNC identifiers, IPI IDs, associated DNA-binding Interpro domains and families, and tissue specificity if any.” Of the 1,987 transcription factors in this supplemental data file, 1,130 were included in the RMA-normalized expression data we had downloaded from GEO and 1,798 had expression information in the normalized GTEx data.

1,120 of these transcription factors had gene expression values in both the RMA-normalized NRG data and the normalized GTEx data (Figure S4A). We evaluated how many of these transcription factors had the same tissue-specific designation in both the NRG supplemental data file and based on our analysis (see Section S4). To do this we created a map between the 38 tissues used in our current GTEx analysis with the 32 tissues analyzed in the NRG paper. In several cases multiple different GTEx tissue subregions (eg the atrial appendage and left ventricle of the heart) were mapped to the same, more general tissue-designation in the NRG data (eg “heart”). We then directly compared the set of transcription factors that were identified as specific to a given tissue in our GTEx analysis, with the set of transcription factors that were identified as specific to that tissue in the NRG analysis.

We find that the overlap between these sets of TFs is nominally statistically significant in most cases ($p < 0.01$ in 14 of the 20 comparisons). However, the actual number of TFs identified as specific to a particular tissue in both the NRG and our GTEx analysis is quite low (Figure S4B). For example, the lung, ovary, and pancreas contained no common tissue-specific TFs between our GTEx designation and the NRG-designation. In addition, when we restrict this analysis to the 478 of these 1,120 TFs that were also included as regulators in our network model, even this nominal significance goes away for many tissues.

To better understand this result, we examined the distribution of expression values in the GTEx data for these 1,120 TFs. A few examples are included in Figure S4C. In some cases, such as for *XBPI* and *TBX3*, the fact that a TF was only identified as specific by NRG and not GTEx appears to be a function of the cutoff we used for defining tissue-specificity (see Section S4). However, we note that relaxing this criterion would have significantly changed the number of TFs we identified as tissue-specific (see Figure S1B) and ultimately would not affect the relatively low level of overlap we see here. In addition, there are many examples where our GTEx analysis clearly identifies tissue-specific signals that are not reflected in the NRG data set (*ZNF106*, *RFX4*, *GATA4*), and also examples where there is no apparent tissue-specific signal for a TF despite it being called so in the NRG data (*EGR1*, *ESR1*). Given that the NRG expression data contains only two samples per tissue, we believe that the tissue-specificity calls for TFs made in our analysis are more reliable.

This low level of overlap in the identified tissue-specific TFs led us to more closely investigate the expression data used in the NRG analysis. Using the RMA-normalized NRG data (and focusing on the 11,363 genes and 1,120 TFs that are common between the NRG and GTEx data sets), we reproduced the plots from Figure 3 in the NRG publication. Consistent with that analysis, we find that in the NRG expression data set transcription factors are expressed at lower levels than non-TFs (compare Figure 3A in [6]) to Figure S4D). We then repeated this same analysis using the GTEx data. To our surprise, the difference in expression between TFs and non-TFs largely disappeared when performing this analysis in the GTEx data. Finally, we repeated this analysis using all 30,333 genes in our GTEx expression data set. This actually resulted in the opposite conclusion as the analysis presented in the NRG paper, with TFs expressed at higher levels than non-TFs.

One advantage of using RNA-sequencing data over microarrays is that sequencing can capture mRNA from many different types of genes and is not limited by the set of probes included on a given array. To better understand whether differences in technology (microarray versus RNA-sequencing) may be influencing the results shown in Figure S4D, we determined the annotations for the 30,333 genes included in our GTEx analysis using Biomart (`dec2013.archive.ensembl.org`). Figure S4E shows the distribution of these annotations across the 11,363 genes that are common between the NRG microarray and the GTEx RNA-seq data, and across the 18,970 genes that are only contained in our GTEx RNA-seq data. It is immediately clear that the microarray genes are almost completely composed of protein-coding genes whereas the genes captured only in the GTEx data contain many types, including antisense, lincRNAs, and pseudogenes. Thus the fact that we see TFs expressed at higher levels than non-TFs when evaluating the full 30,333 genes in the GTEx data is largely a consequence of the fact that all TFs are, by definition, protein-coding genes, and that protein-coding genes are

expressed at higher levels than non-protein-coding genes.

Overall, this analysis highlights the importance of the public availability of data and reproducible research, as we were able to faithfully reproduce many of the results from the NRG paper using their original data. It also highlights the need to revisit previous analyses as new data becomes available. The differences in tissue-specificity and TF-expression based on the NRG analysis and the GTEx data are a perfect demonstration of the opportunity the GTEx data gives us to revisit our understanding of tissue-specificity and gene regulation.

S8. CALCULATING ENRICHMENT OF TISSUE-SPECIFIC EDGES

To quantify the relationship between various tissue-specific edges and nodes, we explicitly evaluated the extent to which tissue-specific edges are more (or less) likely to target tissue-specific genes (or TFs) as compared to chance. For each of the 38 tissues we counted the number of edges called as specific to a tissue (t , see Equation S2), and of a given multiplicity (M , see Equation S3) that also target a gene identified as specific to that tissue (see Equation S4):

$$N^{(t,M)} = \sum_{i,j} \left[(s_{i,j}^{(t)} > 2) \& (m_{ij} == M) \& (s_j^{(t)} > 2) \right]. \quad (S6)$$

We then summed these numbers over all 38 tissues:

$$N^{(M)} = \sum_t N^{(t,M)}. \quad (S7)$$

We also calculated the number of tissue-specific edges of a given multiplicity that one would expect to target tissue-specific genes by chance:

$$\langle N^{(t,M)} \rangle = \frac{1}{N_g} \sum_j (s_j^{(t)} > 2) \sum_{i,j} \left[(s_{i,j}^{(t)} > 2) \& (m_{ij} == M) \right], \quad (S8)$$

$$\langle N^{(M)} \rangle = \sum_t \langle N^{(t,M)} \rangle,$$

where $N_g = 30,243$ (the number of genes in our model). Finally, we defined the enrichment for tissue-specific edges of a given multiplicity targeting tissue-specific genes as:

$$E^{(M)} = \log_2 \frac{Observed}{Expected} = \log_2 \frac{N^{(M)}}{\langle N^{(M)} \rangle}. \quad (S9)$$

We found very high enrichment for tissue-specific edges targeting tissue-specific genes, especially in edges with lower multiplicity values (Figure 3A).

S9. GENE SET ENRICHMENT ON TF TARGETING PROFILES

Gene Set Enrichment Analysis to Quantify the Functions Associated with Tissue-Specific TF-targeting: Although tissue-specific transcription factors are more likely to be associated with tissue-specific network edges than one would expect by chance, we found that this association is much lower than the association between tissue-specific edges and target genes. This led us to the hypothesis that both tissue-specific and non-tissue-specific transcription factors play an important role in mediating tissue-specific biological processes. To test this hypothesis, for each transcription factor (i), we quantified its tissue-specific targeting profile in a given tissue (t) as $s_i^{(t)}$ (see Equation S2). We then ran a pre-ranked Gene Set Enrichment Analysis (GSEA) [10] on the scores in this profile to test for enrichment for Gene Ontology (GO) terms. In total we performed 24,472 GSEA analyses, one for each of the 644 transcription factors included in the network for each of the 38 tissues.

Selection of TFs with Highest and Lowest Expression Enrichment: In order to better understand the relationship between tissue-specific transcription factor expression patterns and their tissue-specific targeting of biological functions, we selected ten transcription factors with the highest expression enrichment based on Equation S4. More specifically, for the analysis presented in Figure 4B in the main text, we selected the ten transcription factors with the highest $s_j^{(Brain\ other)}$ value, and the ten transcription factors for which the absolute value of $s_j^{(Brain\ other)}$ was closest to zero.

Identifying Differentially Targeted Biological Processes and Differentially Targeting TFs for Each Tissue: For each tissue, we identified GO terms that were significantly enriched ($FDR < 0.001$; GSEA Enrichment Score, $ES > 0.65$) for tissue-specific targeting by at least one transcription factor. This allowed us to define 38 sets of differentially targeted biological processes, one for each tissue. For each tissue, we used the corresponding set of differentially targeted GO terms to identify differentially targeting TFs. More specifically, for each tissue we determined the set of TFs that were specifically significantly-enriched ($FDR < 0.001$; GSEA Enrichment

Score, $ES > 0.65$) for differential targeting of at least one differentially targeted biological processes. This allowed us to define 38 sets of differentially targeting TFs, one for each tissue. Interestingly, these TFs were not associated with the sets of differentially expressed (tissue-specific) TFs identified in Section S4 (Figure S5 and Table S5).

Community Structure Analysis to Identify Related Sets of TFs/Tissues and GO terms: To gain a more holistic understanding of the patterns of tissue-specific targeting across all 38 tissues, we combined the GSEA analysis results into a single large matrix that contained the enrichment results across all 24, 472 transcription factor and tissue pairs. This matrix contained all the tested GO terms in the rows, and each of the 24, 472 GSEA analyses in the columns. We selected elements of this matrix that represented highly significant positive enrichment for tissue-specific targeting ($FDR < 0.001$ and $ES > 0.65$), creating a bipartite network where nodes were either GO terms or TF-tissue pairs (the pairs used for the GSEA analyses). We then ran the fast greedy community structure detection algorithm [11] to identify “communities,” or sets of GO terms associated with TF-tissue pairs, in this bipartite network. The benefit of this type of analysis over other clustering approaches, such as hierarchical clustering, is that each “node” is assigned to exactly one community, aiding in our interpretation of these highly complex results. This analysis identified 48 separate communities (Figure 5A and Figure S6), or clusters of GO terms associated with TF-tissue pairs (representing the tissue-specific targeting profile of a particular TF in a particular tissue). All TF-tissue pairs with significant positive enrichment tissue-specific targeting of a particular GO term can be found in Table S5. Characteristics of the 48 communities found by clustering these relationships can be found in Table S6.

Word Clouds to Visualize the Functional Content of Communities: Nine communities had eight or more GO term members. For these communities we summarized their functional content using a free word cloud making program (downloaded from: <http://www.softpedia.com/get/Office-tools/Other-Office-Tools/IBM-Word-Cloud-Generator.shtml>). This program automatically configures the orientation of words in the clouds, but we manually assigned each word a relative size based on that word’s statistical enrichment in the community [12]. Specifically, for a given community, we counted the number of times an individual word appeared across all the GO term members associated with that community (N_{wc}) and then calculated its statistical enrichment in a given community based on the hypergeometric probability:

$$p = \sum_{q=N_{wc}}^{\min[N_w, N_c]} \frac{\binom{N_c}{q} \binom{N_{tot} - N_c}{N_w - i}}{\binom{N_{tot}}{N_w}}, \quad (S10)$$

where N_c is the number of individual words in a community, N_w is the number of times the word appears across all term descriptions and N_{tot} is the total number of words included in all tested GO terms. We then scaled the sizes of the words in the word cloud based on $-\log_{10}(p)$ such that words that have the lowest probability of being in the community by chance are given the largest size and words that are common across many biological functions and that one might expect to be in a community by chance are given a very small size.

Bipartite Network to Visualize the Relationships between Communities and Tissues: Figure 5C was made using JavaScript library D3.js (<http://bl.ocks.org/NPashaP/fcb09e2cddb104e209f457d44f166ca>).

Transcription Factor Enrichment in Communities: Because of the complex structure of the relationships represented between TF/Tissue pairs and the GO Terms in our communities, we identified transcription factors that were significantly enriched in a given functional community by performing a permutation analysis. We began by determining the number of times each transcription factor has a significant GSEA association with each of our 48 functional communities. Then, to determine whether this value was greater than expected by chance, we performed a supervised shuffling of the community labels of the transcription factors. In particular, to perform a community label shuffling that would correctly identify enrichment among transcription factors, we first identified the set of community assignments associated with each tissue, and shuffled these assignments only among the TF/Tissue-pairs for that tissue. This approach allowed us to conserve both the size of the communities and the distribution of tissues within the communities. After performing each shuffling, we counted the number of times each TF had an association with each community in this random assignment. We repeated this shuffling 10,000 times and estimated the significance of enrichment of a transcription factor in each community by determining the percentage of times the counts from the shuffled assignments were greater than the counts from the original assignments.

S10. NETWORK CENTRALITY ESTIMATES OF TISSUE-SPECIFIC GENES

We used the igraph Version 1.0.0 package in R to calculate both the degree (using the `graph.strength()` function) and betweenness centrality (using the `betweenness()` function) of genes in each of the 38 complete, weighted PANDA tissue networks (see Section S2 and Equation S1).

Degree: The degree of a node is defined as the number of edges connected to that node. Because we have weighted graphs, we calculated the degree of a gene in a given tissue (t) by summing up the weights of all edges connected to that gene ($W_j^{(t)}$ see Equation S1). Note that because these are also complete graphs, each gene had exactly 644 edges, one from each transcription factor.

Betweenness: The betweenness of a node is defined as the fraction of non-redundant shortest paths in the network that go through that node. In a weighted network, the shortest path calculation uses edge weights to calculate the cost of traversing each edge. In order to prefer higher edge weights in calculating shortest paths, we used $1/W_{ij}^{(t)}$ (see Equation S1) as the cost for determining the shortest paths. In order to calculate the betweenness centrality, we treated edges as undirected (meaning that an edge exists both from a TF to its target gene and from the target gene to the TF).

S11. NETWORK CENTRALITY OF PANDA'S SEED REGULATORY NETWORK

PANDA builds its predicted regulatory network, in part, by leveraging information from a prior “seed” network constructed by mapping transcription factors to genes based on genome sequence information (see Section S2). We wanted test whether the differences in centrality values that we observed between tissue-specific and non-tissue-specific genes were due to the structure of this input data or if they were identified primarily through PANDA’s message passing network optimization. Therefore, we calculated the degree and betweenness centrality for genes based on the motif scan seed network (see Section S10). We note that this seed network is “un-weighted,” meaning that the edges only take two values: one if the motif for TF i is found in the promoter region of gene j , and zero if it is not.

In the motif prior network, we saw only minimal differences between the centrality of tissue-specific and non-tissue-specific genes, with tissue-specific genes having slightly lower centrality values compared to non-tissue-specific genes (Figure S7). This is consistent with our finding in the main text that tissue-specific genes are generally of low betweenness and only see an increase in their betweenness in their “specific” tissues, and supports our interpretation that tissue specificity is associated with increased centrality in the network as genes gain new non-canonical regulatory paths.

S12. IDENTIFYING GENES WITH EQTLs AND GWAS VARIANT ASSOCIATIONS

To evaluate the potential role of genetic variants in tissue-specific gene regulation, we identified tissue-specific, *cis*-acting expression quantitative trait loci (eQTLs), as described in [13]. Briefly, of the 38 tissues for which we had reconstructed gene regulatory network, 19 contained gene expression samples from at least 150 distinct individuals with imputed genetic data (note that only tissues with at least 200 individuals were presented in [13], but the data were processed the same way). These 19 tissues included adipose subcutaneous, artery aorta, artery tibial, brain other, breast, fibroblast cell line, colon transverse, esophagus mucosa, esophagus muscularis, heart atrial appendage, heart left ventricle, lung, skeletal muscle, tibial nerve, pancreas, skin, stomach, thyroid, and whole blood. For each of these tissues, we identified single-nucleotide polymorphisms (SNPs) that had a minor allele frequency greater than 0.05 across the individuals with associated tissue-specific gene expression data and used Matrix eQTL [14] to quantify the statistical association of the expression of each of the 29, 155 autosomal genes in GTEx with each of these genetic variants. For this analysis we used a *cis*-acting window around the gene of 1 mega-base, and adjusted for sex, age and the three first principal components obtained using the genotyping data. Finally, we determined which genes had at least one significant ($FDR < 0.05$) eQTL association in each tissue.

We also determined which of these QTL-associated genes might be important for disease or other phenotypic traits. To do this, we downloaded the NHGRI-EBI GWAS Catalog (<http://www.ebi.ac.uk/gwas/>; accessed: December 8, 2015). We curated this information, excluding any entries in the catalog for which the variant did not have an associated rsid. Then, we parsed our identified tissue-specific eQTL associations, pruning those that were not with one of these GWAS genetic variants. Finally, we used this information to determine which genes had at least one significant ($FDR < 0.05$) eQTL association with a GWAS-SNP in each tissue. These genes are listed in Table S7.

We compared the genes identified in these analyses with information regarding tissue-specificity. In particular, since we only identified eQTL associations in 19 tissues, we identified a new set of “tissue-specific genes”, which included those that had a multiplicity greater than zero when applying Equation S5 and summing only over the 19 tissues. 5, 256 of the 29, 155 genes were identified as specific to at least one of these 19 tissues.

SUPPLEMENTAL REFERENCES

- [1] GTEx Consortium, *et al.*, “The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans,” *Science* **348**, 6235 (2015).
- [2] S. C. Hicks, O. Kwame, J. N. Paulson, J. Quackenbush, R. A. Irizarry, H. C. Bravo, “Smooth Quantile Normalization,” *Biostatistics*, kxx028 (2017).
- [3] J. N. Paulson, C.-Y. Chen, C. M. Lopes-Ramos, M. L. Kuijjer, J. Platig, A. R. Sonawane, M. Fagny, K. Glass, J. Quackenbush, “Tissue-aware RNA-Seq processing and normalization for heterogeneous and sparse data,” *bioRxiv pre-print* <https://doi.org/10.1101/081802>, (2016).
- [4] K. Glass, C. Huttenhower, J. Quackenbush, G-C Yuan, “Passing messages between biological networks to refine predicted interactions,” *PLoS one* **8**, 5 (2013).
- [5] C. E. Grant, T. L. Bailey, W. S. Noble, “FIMO: scanning for occurrences of a given motif,” *Bioinformatics* **27**, 7 (2011).
- [6] J. M. Vaquerizas, S. K. Kummerfeld, S. A. Teichmann, N. M. Luscombe, “A census of human transcription factors: function, expression and evolution,” *Nature Reviews Genetics* **10**, 4 (2009).
- [7] E. Pierson, D. Koller, A. Battle, S. Mostafavi, GTEx Consortium, *et al.*, “Sharing and specificity of co-expression networks across 35 human tissues,” *PLoS Comput Biol* **11**, 5 (2015).
- [8] C. S. Greene, A. Krishnan, A. K. Wong, E. Ricciotti, R. A. Zelaya, D. S. Himmelstein, R. Zhang, B. M. Hartmann, E. Zaslavsky, S. C. Sealfon, *et al.*, “Understanding multicellular function and disease with human tissue-specific networks,” *Nature Genetics* **47**, 569 (2015).
- [9] M. Dai, P. Wang, A. D. Boyd, G. Kostov, B. Athey, E. G. Jones, W. E. Bunney, R. M. Myers, T. P. Speed, H. Akil, *et al.*, “Evolving gene-transcript definitions significantly alter the interpretation of GeneChip data,” *Nucleic acids research* **33**, 20 (2005).

- [10] A. Subramanian, P. Tamayo, C. K. Mootha, S. Mukherjee, B. L. Ebert, M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander, *et al*, "Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles," *Proceedings of the National Academy of Sciences* **102**, 43 (2005).
- [11] A. Clauset, M. E. J. Newman, C. Moore, "Finding community structure in very large networks," *Physical review E* **70**, 6 (2004).
- [12] K. Glass, M. Girvan, "Finding new order in biological functions from the network structure of gene annotations," *PLoS Comput Biol* **11**, 11 (2015).
- [13] M. Fagny, J. N. Paulson, M. L. Kuijjer, A. R. Sonawane, C-Y. Chen, C. M. Lopes-Ramos, K. Glass, J. Quackenbush, J. Platig, "A network-based approach to eQTL interpretation and SNP functional characterization," *Proceedings of the National Academy of Sciences*, doi: 10.1073/pnas.1707375114, *ePub ahead of print*.
- [14] Shabalin, "Matrix eQTL: ultra fast eQTL analysis via large matrix operations," *Bioinformatics* **28**, 1353 (2012).