

## ARTICLE OPEN



# Multilayer modelling of the human transcriptome and biological mechanisms of complex diseases and traits

Tiago Azevedo<sup>1,7</sup>, Giovanna Maria Dimitri<sup>1,2,3,7</sup>, Pietro Lió<sup>1,2</sup>✉ and Eric R. Gamazon<sup>2,4,5,6</sup>✉

Here, we performed a comprehensive intra-tissue and inter-tissue multilayer network analysis of the human transcriptome. We generated an atlas of communities in gene co-expression networks in 49 tissues (GTEx v8), evaluated their tissue specificity, and investigated their methodological implications. UMAP embeddings of gene expression from the communities (representing nearly 18% of all genes) robustly identified biologically-meaningful clusters. Notably, new gene expression data can be embedded into our algorithmically derived models to accelerate discoveries in high-dimensional molecular datasets and downstream diagnostic or prognostic applications. We demonstrate the generalisability of our approach through systematic testing in external genomic and transcriptomic datasets. Methodologically, prioritisation of the communities in a transcriptome-wide association study of the biomarker C-reactive protein (CRP) in 361,194 individuals in the UK Biobank identified genetically-determined expression changes associated with CRP and led to considerably improved performance. Furthermore, a deep learning framework applied to the communities in nearly 11,000 tumors profiled by The Cancer Genome Atlas across 33 different cancer types learned biologically-meaningful latent spaces, representing metastasis ( $p < 2.2 \times 10^{-16}$ ) and stemness ( $p < 2.2 \times 10^{-16}$ ). Our study provides a rich genomic resource to catalyse research into inter-tissue regulatory mechanisms, and their downstream consequences on human disease.

*npj Systems Biology and Applications* (2021)7:24; <https://doi.org/10.1038/s41540-021-00186-6>

## INTRODUCTION

The modern science of networks has contributed to notable advances in a range of disciplines, facilitating complex representations of biological, social, and technological systems<sup>1</sup>. A key aspect of such systems is the existence of community structures, wherein groups of nodes are organized into dense internal connections with sparser connections between groups. Community structure detection in genome-wide gene expression data may enable detection of regulatory relationships between regulators (e.g., transcription factors or microRNAs), and their targets and capture novel tissue biology otherwise difficult to reach. Furthermore, it offers opportunities for data-driven discovery and functional annotation of biological pathways.

We hypothesize that community structure is an important organizing principle of the human transcriptome, with critical implications for biological discovery and clinical applications. Co-expression networks, in fact, encode functionally relevant relationships between genes, including gene interactions and coordinated transcriptional regulation<sup>2</sup>, and provide an approach to elucidating the molecular basis of disease traits<sup>3</sup>. Therefore, reconstructing communities of genes in the transcriptome may uncover novel relationships between genes, facilitate insights into regulatory processes, and improve the mapping of the human diseasome.

In this work, we develop a model of the human transcriptome as a multilayer network, and perform a comprehensive analysis of the communities obtained with this modeling in order to further our understanding of its wiring diagram and facilitate research into improved disease diagnosis and profiling. We conduct a systematic analysis of the tissue or cell-type specificity of the

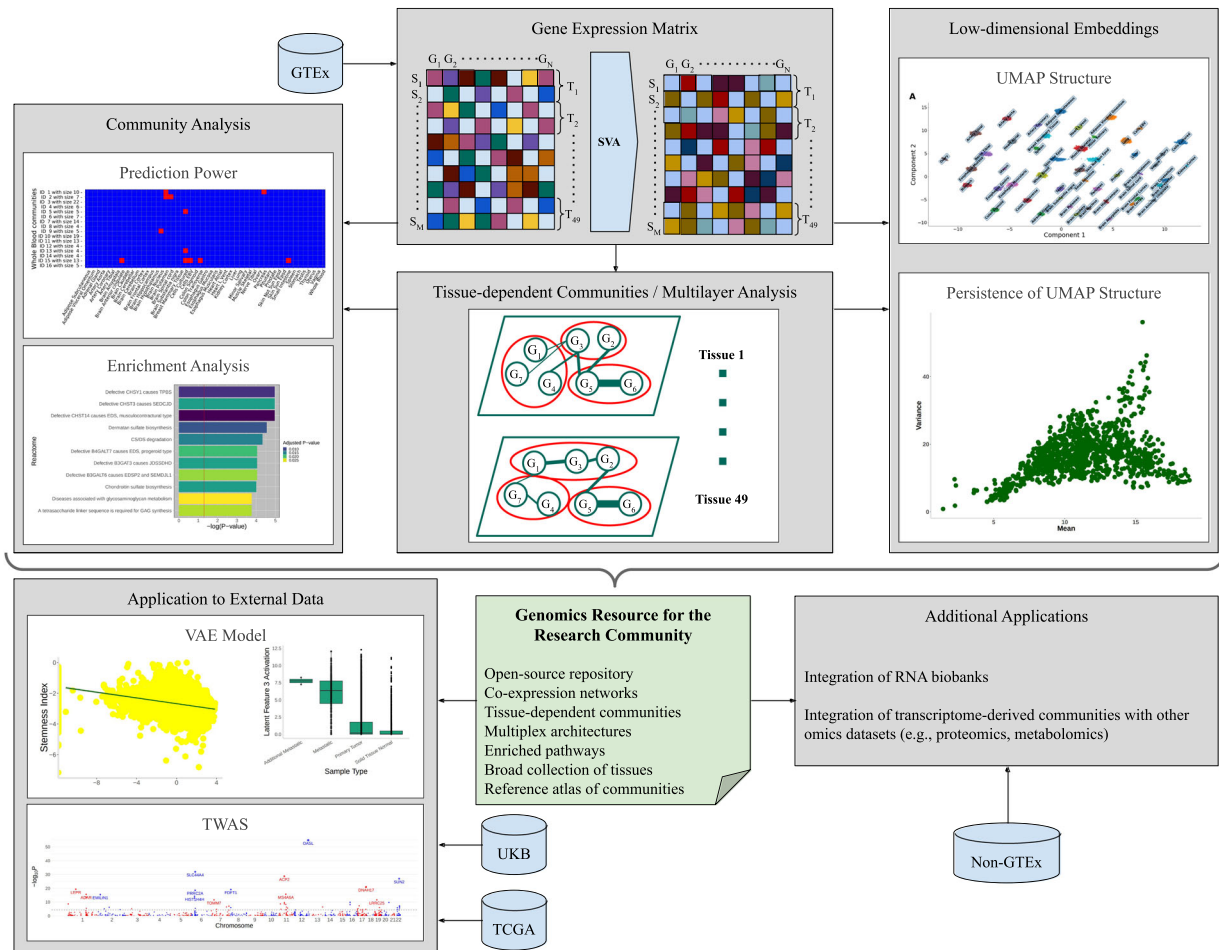
communities in the transcriptome in order to gain insights into gene function in the genome, and enhance our ability to identify disease-associated genes. Our study represents an effort to fill an important gap in our understanding of the role of gene expression in complex traits, i.e., how a gene's phenotypic consequence on disease or trait<sup>4</sup> is mediated by its membership in tissue-specific biological modules as molecular substrates. Methodologically, we demonstrate an approach to integrating the communities into transcriptome-wide association studies (TWAS)<sup>5-7</sup>, and a deep neural network methodology for generating biologically-meaningful latent representations of gene expression<sup>8,9</sup>. Finally, the inter-tissue analysis of the transcriptome holds promise for identifying novel regulatory mechanisms, enhancing our understanding of trait variation and pleiotropy, and opening up new possibilities for translational applications.

## RESULTS

### Study design

Here, we provide a brief overview of our study design. We performed a comprehensive intra-tissue and inter-tissue network analysis of the human transcriptome. We leveraged the GTEx v8 dataset to generate an atlas of communities in co-expression networks in 49 human tissues. Furthermore, we investigated the methodological implications of the communities derived from gene expression. Figure 1 is a schematic of our analytic workflow.

<sup>1</sup>Department of Computer Science and Technology, University of Cambridge, Cambridge, UK. <sup>2</sup>Clare Hall, University of Cambridge, Cambridge, UK. <sup>3</sup>Department of Engineering, University of Siena, Siena, Italy. <sup>4</sup>Vanderbilt Genetics Institute, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>5</sup>Data Science Institute, Vanderbilt University, Nashville, TN, USA. <sup>6</sup>MRC Epidemiology Unit, University of Cambridge, Cambridge, UK. <sup>7</sup>These authors contributed equally: Tiago Azevedo, Giovanna Maria Dimitri. ✉email: pl219@cam.ac.uk; ericgamazon@gmail.com



**Fig. 1 Study design.** We leverage the matrix of gene expression in 49 GTEx tissues. Surrogate variable analysis (SVA) is applied to the high-dimensional dataset to adjust for unknown or unmodelled confounders. Tissue-dependent communities are generated, analysed, tested for enrichment for known biological processes, and exploited towards identification of new functional gene sets. Uniform Manifold Approximation and Projection (UMAP) embeddings of gene expression data defined by the communities, and the persistence of the global structure, are evaluated to identify biologically-meaningful clusters. Notably, new datasets can be embedded into the derived models to facilitate additional discoveries. Potential external applications of the resource of gene expression communities are varied. Here, we implement two community-based applications, including a deep learning (variational autoencoder) model and transcriptome-wide association studies (TWAS) using the TCGA and the UK Biobank datasets, respectively.

### Spurious co-expression and confounding due to unmodelled factors

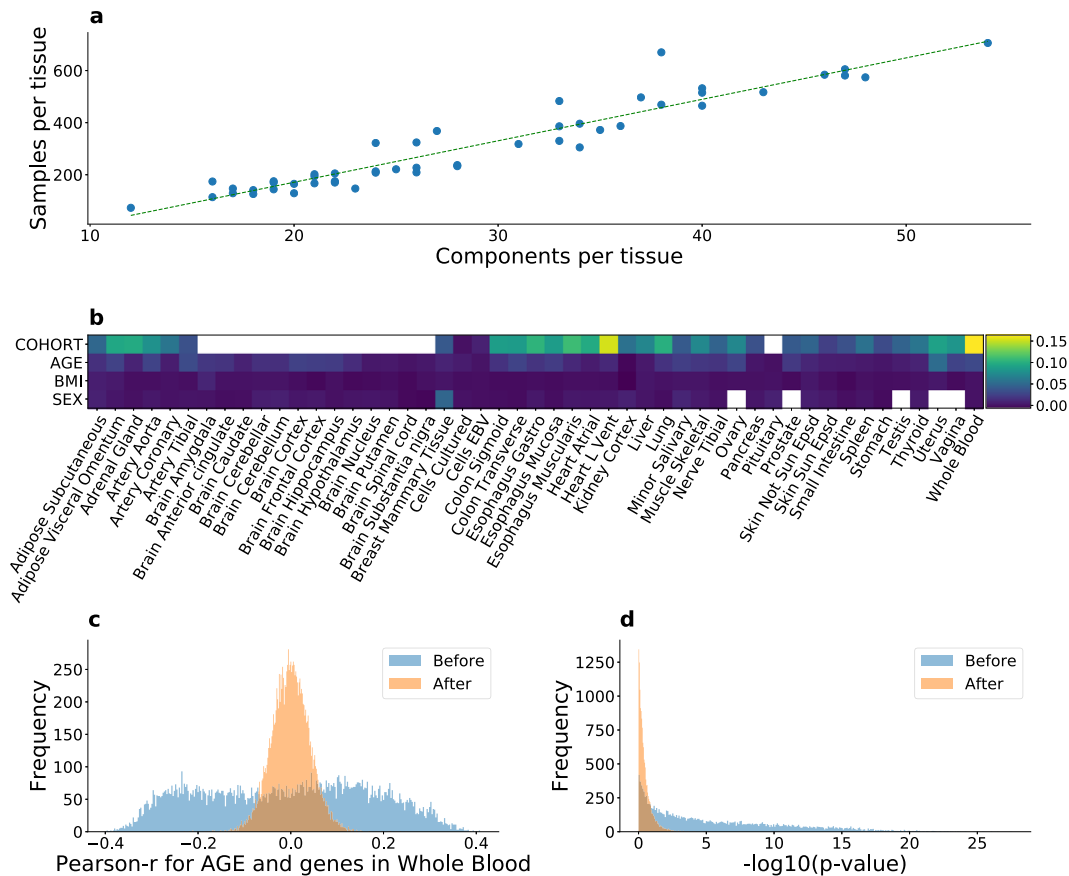
Disambiguating true co-expression from artefacts is an important concern in the presence of hidden variables. We therefore applied *sva* analysis to investigate unmodelled and unmeasured sources of expression heterogeneity<sup>10</sup>. The number of factors or components identified by this analysis was significantly correlated ( $r \approx 0.95$ ,  $p \approx 5.4 \times 10^{-26}$ ) with the number of distinct samples across tissues (see Fig. 2a). Notably, the greater number of such surrogate variables that we regressed out for tissues with larger sample sizes recapitulates the approach used by the GTEx Consortium of using more inferred factors for tissues with larger sample sizes, in order to optimize the number of eGenes from the eQTL analysis<sup>11</sup>. Specifically, the GTEx Consortium uses PEER, a related adjustment method, with 15 factors for tissue sample size  $N < 150$  and up to 60 factors for  $N \geq 150$ .

We then quantified the impact of confound correction (see Fig. 2b) in co-expression analysis. The distribution of Pearson correlation values has more mass closer to zero with less variance after correction, suggesting that unmodelled factors may induce spurious (or artificially inflate) correlations in gene expression. The

effect of unmodelled factors is further illustrated in Fig. 2b–d, where the distribution of correlation values for the covariate *Age* is shown for whole blood. Before correction, those values are spread between around  $-0.4$  and  $0.4$ , whereas after correction the corresponding values move towards the center (zero) and become less dispersed. Notably, the variable *Cohort* (with possible values being *Postmortem* and *Organ Donor* in available tissues, except for some which also have *Surgical* values) seems to have undergone the largest change in the correction process. This suggests that estimation of cohort effect on gene expression can be substantially improved by accounting for unmodelled factors.

### Atlas of communities across human tissues

For each tissue, we identified communities in the co-expression networks, using the Louvain algorithm (see “Methods” section), to develop an atlas across human tissues. On average, a tissue was found to have 108 communities (standard deviation [SD] = 31) (see Fig. 3). We observed the highest number of communities ( $n = 251$ ) in “Kidney cortex” and the lowest number ( $n = 73$ ) in “Muscle skeletal”. The nonsolid tissues, consisting of “Cells EBV” and “Whole blood”, have the highest number of genes (i.e., at least



**Fig. 2 Confounding due to unmodelled factors.** **a** Relationship between the number of inferred factors and tissue sample size. Fitted line ( $r \approx 0.95$ ,  $p \approx 5.4 \times 10^{-26}$ ) corresponds to a linear least-squares regression. The two-sided  $p$ -value is based on the null hypothesis that the slope is zero, using the Wald Test with  $t$ -distribution for the test statistic. **b** The difference in the variance of the distribution of Pearson correlation values for each tissue over all genes, before and after correction. Empty cells correspond to tissues in which only one value of the confound is available. The "Cohort" variable undergoes the most substantial change after the correction across all tissues. **c** Distribution of Pearson correlation between the expression of a gene in whole blood and age, before and after correction. After the correction, the correlation values move towards zero and show considerably less dispersion. **d** The  $p$ -value distribution from panel **c**'s, in logarithmic space. The enrichment for significant (low)  $p$ -values is greatly attenuated after the correction, suggesting that unmeasured variables can induce spuriously significant correlations.

4300 for each) that belong to a community. The size of a community varies considerably within each tissue and its distribution differs across tissues (see Supplementary Table 1 and Supplementary Fig. 1 for the distribution in all tissues). The brain tissues show significantly higher variability (median SD = 9.9, Mann–Whitney  $U$ -test  $p = 1.55 \times 10^{-4}$ ) than non-brain tissues (median SD = 5.18). Thus, tissues and tissue classes may differ in the overall topology of the communities in co-expression networks, which likely contains considerable tissue information.

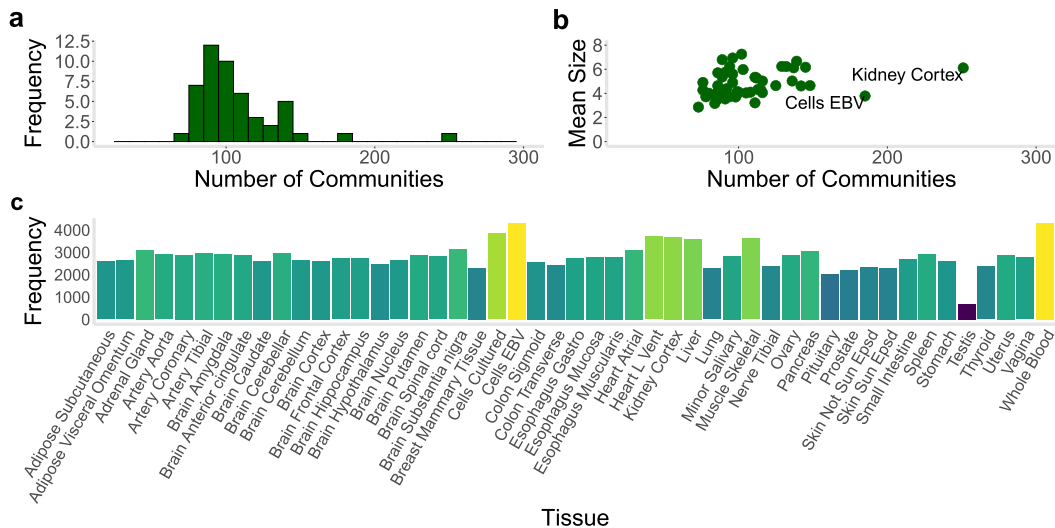
We noticed that after removing the weaker correlations ( $-0.80 < z_{ij} < 0.80$ ), most of the subnetworks were already highly segregated from the rest of the entire network, indicating that the Louvain communities could be completely formed by just this removal process. In order to evaluate the segregation of such communities, we calculated the number of connections coming out of communities of each size. We found that for every tissue the mode was zero, and the maximum number was never over 17. Given the thousands of genes in each tissue's co-expression network, the observed maximum number of connections between different communities (i.e., at most 17) illustrates how strong the segregation is prior to the application of the Louvain community analysis.

More information on these communities is available on github repository (see "Code availability" section): notebook *09\_community\_info*.

### UMAP of community-defined gene expression manifold reveals tissue clusters

To generate a lower dimensional representation of the original transcriptome dataset, we performed Uniform Manifold Approximation and Projection (UMAP)<sup>12</sup> (see "Methods" section). Nearly 18% of the genes belong to a community in at least one tissue. Notably, gene expression from this subset was able to recover the tissue clusters (see Fig. 4a) as fully as the complete set of genes analysed here (see Supplementary Fig. 2).

Drawing conclusions about relationships between clusters (tissues) from UMAP and similar approaches must be done with caution due to some known caveats<sup>13</sup> (see next section for more details). However, starting from known relationships between tissues, we found that the subset of community-based genes yielded biologically consistent embeddings from UMAP. Indeed, the clustering of related tissues (based on organ membership), such as the 13 brains regions, or the clustering of other related tissues (based on shared function), such as the hypothalamus-pituitary complex (which controls the endocrine system<sup>14</sup>), could be observed for the genes that belong to communities. Taken together, these results show that gene expression from the identified communities encodes sufficient information to distinguish the various tissues in a biologically-meaningful low-



**Fig. 3 Summary statistics on identified communities.** **a** Histogram shows the distribution of community count in the various tissues (mean = 108, SD = 31). **b** The scatter plot displays the community count and mean community size for each tissue, showing a significant correlation (Spearman  $\rho = 0.39$ ,  $p = 0.006$ ). The highest number of communities was observed in "Kidney Cortex" ( $n = 251$ ). **c** Plot provides the number of genes that belong to a community in each tissue. The nonsolid tissues, "Cells EBV" and "Whole Blood", show the highest number of genes with membership in a community.

dimensional representation. We note, however, that not all sets of genes with correlated expression produce the distinct separation of tissues observed for the set of genes that belong to the communities (see below).

In theory, additional clusters may be present at different scales, such as within a tissue. Therefore, we performed UMAP analysis on the single-tissue "Whole Blood" to test for the presence of additional clusters. Notably, no well-defined clustering was observed with respect to cohort (Supplementary Fig. 3), BMI (Supplementary Fig. 4), and the other covariates, indicating that the *sva* analysis was successful in removing potential confounders (see Fig. 2b).

External transcriptome data can be embedded into the trained model generated from the GTEx communities. Indeed, using TCGA data for 33 cancer types, we found that the embeddings into the learned space recapitulate recent findings on cancer-testis (CT) genes (see Supplementary Material). In addition, UMAP representations of the genes that belong to a GTEx-derived community recovered the cancer types (see Supplementary Fig. 5) in (external) TCGA data, showing the cross-study relevance of the model.

### Persistence of the UMAP embeddings

We developed an approach to quantify the conservation and variability of the UMAP global structure and estimate the sampling distribution of the local structure, i.e., the distance  $d(i, j)$  for a given pair of tissues  $i$  and  $j$  (see "Methods" section). Using 500 bootstrapped manifolds, we found that on average related tissues tended to cluster closely together (see Fig. 4b). Examples of such clusters are the 13 brain regions, the colonic and esophageal tissues, and various artery tissues. Supplementary Fig. 6 shows the relationship between the average distance between tissue clusters and the variance in the distance, showing a significant positive correlation (Spearman  $\rho \approx 0.38$ ,  $p < 2.2 \times 10^{-16}$ ). Reassuringly, the tissue pairs ("Brain Cerebellum", "Brain Cerebellar") and ("Skin Not Sun Epsd", "Skin Sun Epsd") had the lowest average distance between clusters among all tissue pairs; the first pair consists of known duplicates of a brain region in the GTEx data<sup>15</sup> and is thus expected to cluster together. Among the tissue pairs with the highest average distance, "Adipose Subcutaneous" had an average distance greater than 17 with each of the colonic tissues ("Colon Sigmoid" and "Colon Transverse"), and a low variance comparable

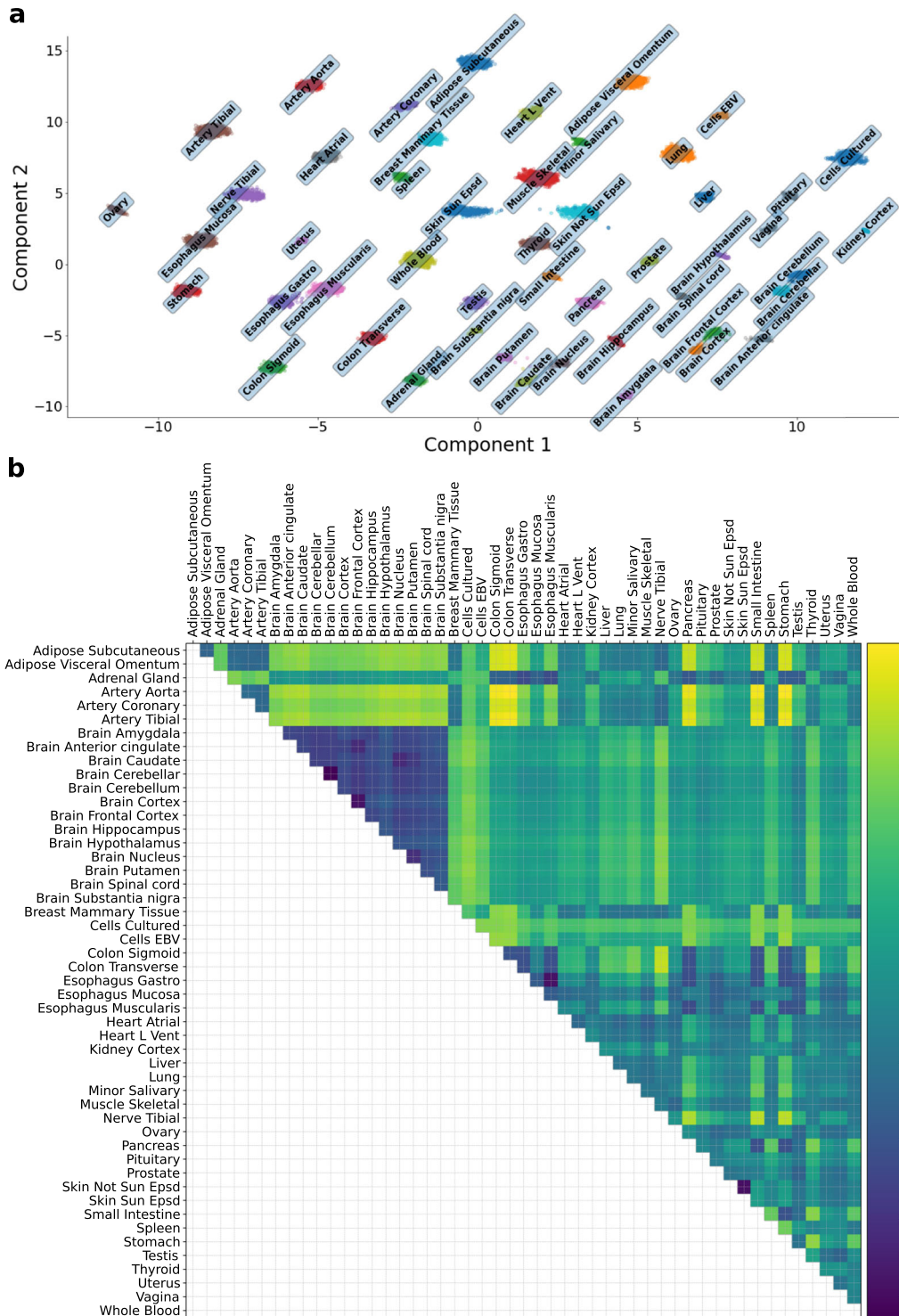
to tissue pairs with some of the smallest average distance. Additional global patterns can be easily observed. For example, the relationships of related tissues (e.g., "Skin Sun Epsd" and "Skin Not Sun Epsd" with  $C_{(i_0, i_1)} \approx 0.62$ ,  $p = 3.4 \times 10^{-5}$ ) to the remaining tissues were found to be strongly preserved, using our clustering conservation coefficient (see "Methods" section).

We also quantified the conservation and variability of the UMAP global structure using the TCGA data from 33 cancer types. The application of the GTEx communities (with only 18% of all analysed genes) in the TCGA data generated biologically consistent UMAP clustering of the cancer types (see Supplementary Fig. 7).

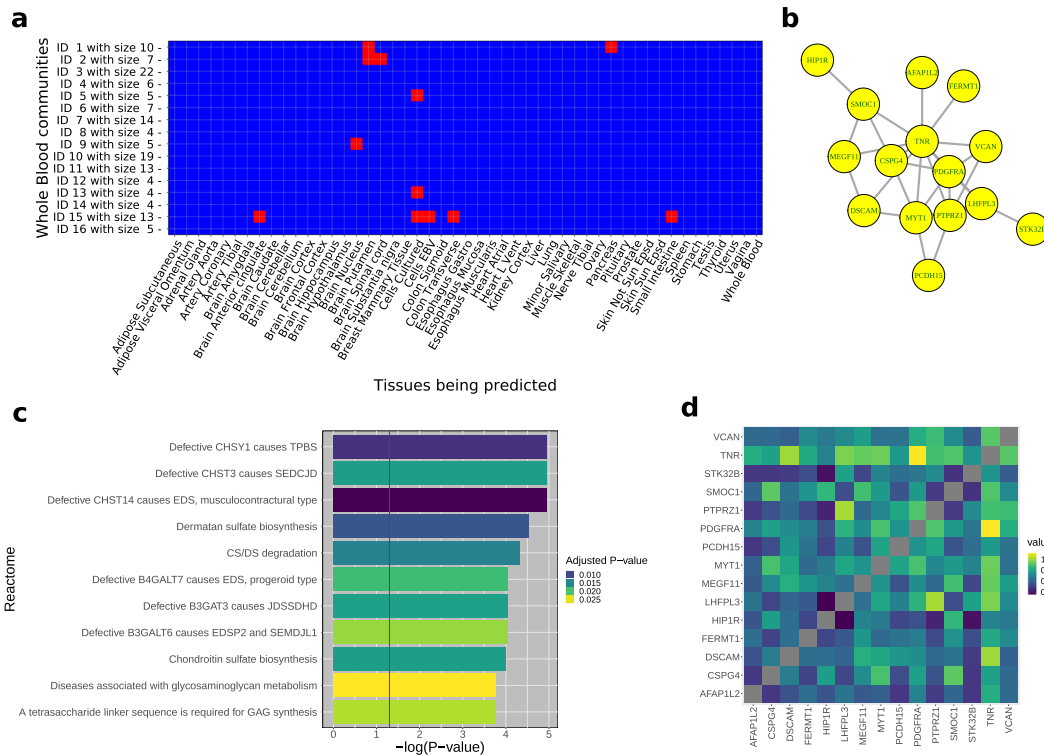
### Prediction of tissues by communities

We then tested individual communities for their ability to predict a tissue. By definition, we consider that a set of genes can predict a tissue when the average  $F_1$  score is above 0.80 (see "Methods" section). Some broad patterns are noteworthy. Most of the communities from "Whole Blood" do not have prediction power for the other tissues (Fig. 5a) partly due to the stringency of our  $F_1$  threshold, which is likely to produce false negatives. This observation indicates that the member genes in each such community from the source tissue ("Whole Blood") cannot "separate" the test tissue (say, "Lung") from the remaining tissues possibly, due to lack of tissue specificity of the gene expression profile of the community. However, a community of only five genes can predict the brain region nucleus accumbens (basal ganglia). To perform the classification, we used a linear classifier, and the so-called "kernel trick" may work in the non-linearly separable gene expression profiles, though perhaps at the expense of biological interpretability. For these communities, the member genes, collectively, are "differentially expressed" between the test brain region and the remaining tissues. Thus, although the genes are present in "Whole Blood" (as a community), the expression profile in the test brain region is substantially different or tissue-specific. "Cells cultured fibroblasts" is the tissue which can be predicted by the largest number of "Whole Blood" communities (three) and, consistently, the largest number from the other tissues.

We note that, consistent with our observations for the communities, 197 Reactome pathways are not sufficient to predict



**Fig. 4 Lower-dimensional UMAP representation of the transcriptome data restricted to the communities and conservation of global structure.** **a** UMAP generates embedded structures through a low-dimensional projection of the submatrix consisting of only the genes that belong to a community in at least one tissue ( $n = 3259$ ). This subset of genes (17.7% of total) contains sufficient information to recover the tissue clusters. In addition, known relationships between tissues, based on organ membership and, separately, on shared function, are reflected in the UMAP projection. **b** Using bootstrapped manifolds (see “Methods” section), we estimated the persistence of the global structure and pairwise relationships across tissue clusters. Here, we show the upper-triangular matrix of the average pairwise distances across the bootstrapped manifolds. We found consistent clustering of known related tissues, including the 13 brain regions, the colonic and esophageal tissues, and various artery tissues. Additional patterns were observed. For example, as reflected in the heatmap, we found a highly correlated relationship, i.e., high “clustering conservation coefficient” ( $C_{(i_0, i_1)} \approx 0.62$ ,  $p = 3.4 \times 10^{-5}$ ) (see “Methods” section), of the two skin tissues to all the other tissues.



**Fig. 5 Communities and their properties.** We developed tools (available on github) to query a community for its characterization. **a** Prediction power of "Whole Blood" communities, in  $F_1$  scores thresholded over 0.8. Most communities in "Whole Blood" do not have prediction power for the remaining tissues. Notable exceptions include a 13-member community, which can predict multiple tissues, including "Brain Anterior cingulate", "Small Intestine", and "Colon Transverse". **b** A 15-member community in the hippocampus is shown here as an example. An edge indicates  $A_{ij} > 0.80$  for genes  $i$  and  $j$ . The gene *TNR*, which is expressed primarily in the central nervous system and involved in its development, is connected via an edge to twelve member genes while *HIP1R* is connected to only one. **c** Enrichment analysis was performed on all communities to identify known biological processes. For example, the hippocampal community in **b** was found to be significantly enriched for Reactome pathways.  $p$ -value refers to raw  $p$ -value. Red line corresponds to the raw  $p < 0.05$  threshold. Color gradient reflects the adjusted  $p$ -value. All Reactome pathways shown meet adjusted  $p < 0.05$ . **d** Heatmap displays the correlation values for the member genes of the community in **b**.

any tissue (available on github: output *output\_06\_02*), while 164 are tissue-specific (i.e., can predict only one tissue) (see details in Supplementary Material).

### Enrichment of communities for known biological processes

We quantified the extent to which the communities in the various tissues reflect current biological knowledge (as encoded in the Reactome pathways). We identified 114 communities (8.28% of all the communities with more than three member genes) enriched for some Reactome pathway (i.e., at an adjusted  $p < 0.05$  for level of enrichment), thus contributing in complex ways to multiple biomolecular processes. "Whole Blood" was the only tissue without any community enriched for known pathways, and the "Esophagus Mucosa" was the tissue with the most communities enriched for known pathways, with a total of five communities. Since the entire set of communities could fully recover all tissues as clusters in the UMAP embeddings, these results suggest that the remainder of the communities are likely to capture previously inaccessible and novel tissue biology.

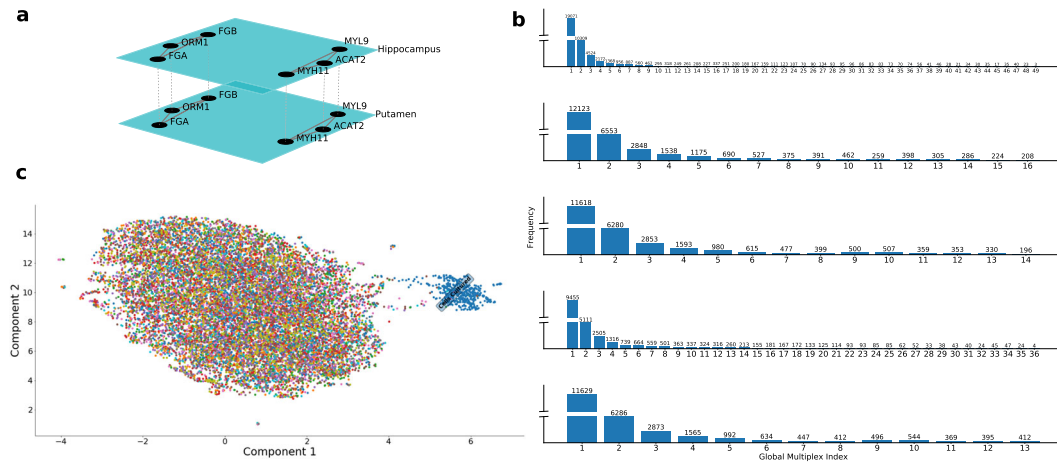
Notably, our analysis may uncover the role of these communities in human diseases. For example, a community of 15 genes in the "Brain Hippocampus" showed a significant enrichment for diseases associated with glycosaminoglycan metabolism (adjusted  $p$ -value = 0.026; see Fig. 5). Glycosaminoglycans, which are major extracellular matrix components whose interactions with tissue effectors can alter tissue integrity, have been shown to play a role in brain development<sup>16</sup>, modulating neurite outgrowth and participating in synaptogenesis. Alterations of glycosaminoglycan

structures from Alzheimer's disease hippocampus have been implicated in impaired tissue homeostasis in the Alzheimer's disease brain<sup>17</sup>.

### Multiplex analysis of the transcriptome

We analysed five multiplex networks to model the various tissue interactions, of clear biological interest, in the GTEx dataset (see "Methods" section). For each multiplex architecture, only the specific component tissues were used to construct the multiplex network, and consequently we calculated the global community index for each multiplex architecture separately. The five architectures analysed were:

- All tissues: Each layer represents one of the 49 tissues analysed. This architecture allows us to investigate gene communities that are shared across all tissues, with potentially universal function.
- Brain and gastrointestinal tissues: The 16 layers correspond to the brain tissues and three gastrointestinal tissues. This architecture may provide insights into the gut-brain axis, which has attracted recent attention in the literature, such as in studies of neuropsychiatric processes and of the interaction between the CNS and the enteric nervous system in neurological disorders<sup>18,19</sup>.
- Brain tissues and whole blood: This multiplex model consists of the 14 layers corresponding to these tissues. This architecture allows us to study brain-derived communities for which the easily accessible whole blood can serve as a



**Fig. 6 Multiplex analysis.** **a** An example of a multiplex network in "Brain Hippocampus" and "Brain Putamen" (basal ganglia). Each layer denotes a tissue, and nodes are the genes, which are connected via the interlayer connections. In the multiplex example, we see the presence of two communities. These two communities are indeed part of the "Brain Tissues" multiplex architecture, present in all 13 layers (brain regions). All genes in the community {*FGA*, *ORM1*, *FGB*} have been implicated as biomarker and therapeutic targets for intracerebral hemorrhage. **b** Histograms show the empirical distribution of the global multiplexity index for each multiplex architecture (with positive index). The index quantifies how many times two genes belong to the same communities across layers. The proportion at each value  $k$  of the index is an estimate of the  $\pi_k$  (see "Methods" section). The maximum value corresponds to the number of layers or tissues of the multiplex network. Histograms, from the top, correspond to: all tissues, brain and gastrointestinal tissues, brain tissues and whole blood, non-brain tissues, brain tissues. **c** We performed UMAP on the subset of tissues that exist across all layers of the central nervous system ("Brain Tissues") multiplex. This set does not yield complete clustering of tissues.

proxy tissue.

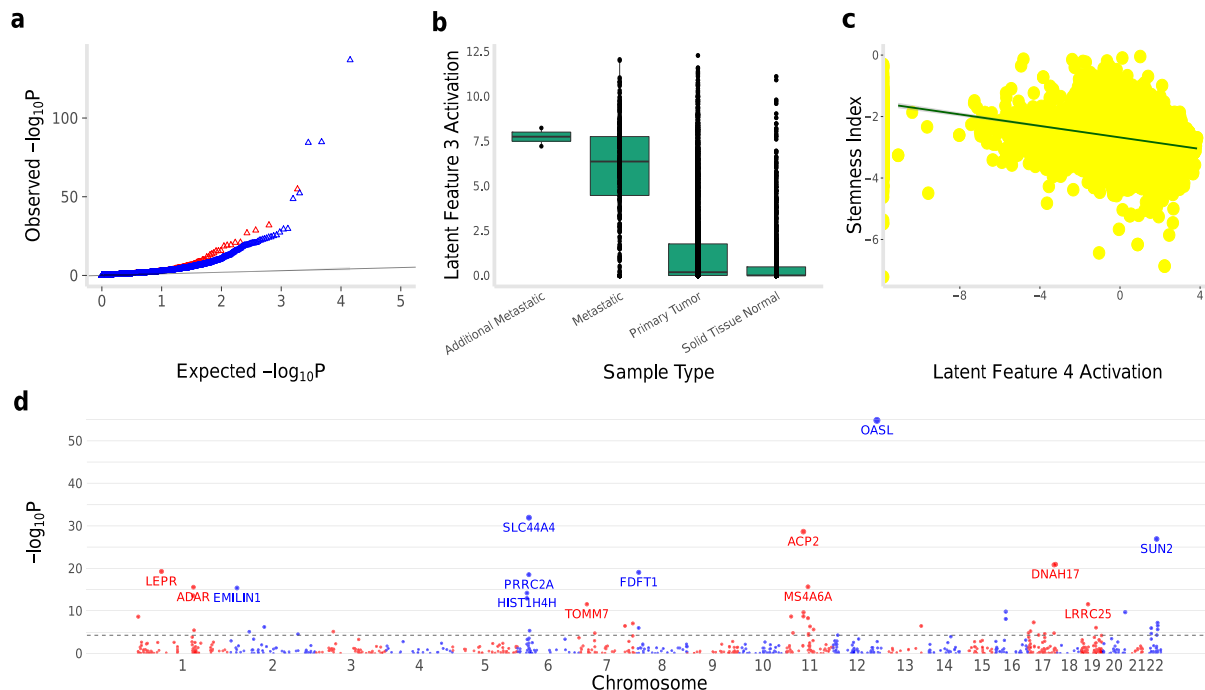
- Non-brain tissues: The 36 layers consist of all tissues outside the brain. This architecture may stimulate investigations into developmental and pathophysiological processes outside the CNS.
- Brain tissues: The 13 layers correspond to the various brain regions. This architecture facilitates identification of communities that may play a functional role throughout the central nervous system (CNS).

Multiplex analysis provides an inter-tissue framework for the analysis of high-dimensional molecular traits, such as gene expression. The global multiplexity matrix was obtained for each of the five proposed architectures. We extracted from the global multiplexity matrices the groups of genes with the maximal global multiplexity index in the five architectures, i.e., the groups of genes that share a value of 49, 16, 14, 36, and 13 respectively, equal to the number of layers (tissues) in the respective architectures. Among these groups of genes with the highest global multiplexity index, we obtained the sub-clusters for each architecture, identifying the groups of genes that always appear in the same community across the various layers. Revealing the shared community structure across the layers improves our understanding of the functional and disease consequences of the clusters of genes. We investigated the biological pathways (Reactome) in which such subgroups were involved for each architecture. Our goal was to test the communities for enrichment for known biological pathways, and therefore quantify the degree to which the communities capture current understanding of biological processes as encoded in the knowledge base.

We illustrate this approach here. In Fig. 6a, we show an example of a multiplex network. In this case, the multiplex network is constructed from data in two brain tissues: "Brain Hippocampus" and "Brain Putamen" (basal ganglia). Each layer represents a tissue, and nodes are labeled as the genes to which they correspond. In the example, we see the presence of two communities formed by groups of genes randomly drawn from the full set: {*FGA*, *ORM1*, and *FGB*} and {*MYH11*, *ACAT2*, and *MYL9*}, with the intralayer and interlayer connections shown. We note that, based on our analysis, these two communities are indeed part of the "Brain Tissues"

multiplex architecture, present in all 13 component layers (brain regions). Notably, all three genes that belong to the first community have been previously implicated as biomarker and therapeutic target candidates for intracerebral hemorrhage<sup>20</sup>. This observation is interesting given our Reactome pathway analysis results; the community is enriched for "Common Pathway of Fibrin Clot Formation" (adjusted  $p = 8.8 \times 10^{-4}$ ) and "Formation of Fibrin Clot (Clotting Cascade)" (adjusted  $p = 1.7 \times 10^{-3}$ ), indicating the genes' involvement in coagulation. All three members of the second community are myelin-associated genes<sup>21</sup>. We found a 17-member community in the "Brain Tissues" multiplex that is significantly enriched for the nonsense-mediated decay (NMD) pathway (adjusted  $p = 1.01 \times 10^{-37}$ ), which is known to be a critical modulator of neural development and function<sup>22</sup>. The pathway accelerates the degradation of mRNAs with premature termination codons, limiting the expression of the truncated proteins with potentially deleterious effects. The community's presence in all brain regions underscores its crucial protective function throughout the central nervous system.

The multiplex analysis we performed can also be used to investigate the relationship between two distinct systems. Here we illustrate this using the CNS and the gastrointestinal system, possibly reflecting a coordinated transcriptional regulatory mechanism between the CNS and the enteric nervous system (ENS). The ENS is a large part of the autonomic nervous system that can control gastrointestinal behavior<sup>23</sup>. We found a 14-member community in the "Brain and Gastrointestinal Tissues" multiplex, whose presence in all 16 layers suggests a strong interaction between (in addition to shared function across) the CNS and ENS. Consistent with this hypothesis, the community was found to be significantly enriched for the "metabolism of vitamins and cofactors" (adjusted  $p = 6.5 \times 10^{-7}$ ), which has been shown to be responsible for altered functioning of the CNS and ENS<sup>24</sup>. Although the involvement of the individual member genes in this pathway is known, the finding that the genes are organized as a community structure, within co-expression networks, which persists across the entire 16 layers of the various brain regions and the gastrointestinal tissues examined here was made possible by the unique transcriptome dataset.



**Fig. 7 Integrating communities into TWAS and Variational Autoencoder.** Leveraging the communities in genomic studies of disease may improve identification of disease-associated genes and enable unique biological insights. **a** We performed a transcriptome-wide association study (TWAS) of C-reactive Protein (CRP) in 361,194 UK Biobank individuals using PrediXcan. CRP is a biomarker for a wide range of complex diseases. The set of genes in the communities (red) displayed a greater departure from the null expectation (i.e., greater enrichment for significant associations) than the complement set of genes (blue), as shown by the leftward shift in the Q-Q plot. **b** We implemented a variational autoencoder (VAE) to leverage the power of neural networks in high-dimensional molecular data. Using 11,060 samples across 33 cancer types from TCGA, the VAE learned biologically-meaningful latent spaces from the communities. For example, latent feature three separates metastatic tumors from primary tumors and normal solid tissues (Mann–Whitney  $U$ -test  $p < 2.2 \times 10^{-16}$  for each comparison), representing a potential mechanism for metastasis. Median line and scatter points are shown. Box edges show interquartile range. **c** The VAE model learned a stemness index significantly associated ( $p < 2.2 \times 10^{-16}$ ) with a DNA-methylation-based index. Least-squares regression line is shown. **d** Manhattan plot shows TWAS associations with CRP for the genes in the communities (dashed line is Bonferroni-adjusted  $p < 0.05$ ).

The empirical distribution of the global multiplexity index is presented in Fig. 6b for each of the five architectures. The maximal global multiplexity index in the five architectures represents the groups of genes that share a value of 49, 16, 14, 36, and 13 respectively, equal to the number of layers (i.e., tissues) in the respective architectures. These genes appear in the same community across all layers of the respective architectures.

For comparison with the UMAP embeddings of the set of all communities, we performed similar analyses in the various multiplex networks. For example, we tested whether the complete tissue clustering could be observed using just the subset of communities that exist across all layers of the central nervous system multiplex. We discovered a different clustering pattern, with cultured fibroblasts clustering separately from the rest of the tissues, which no longer show well-defined clustering (Fig. 6c). This finding suggests the presence of a hierarchy of clusters in the transcriptome at increasingly finer scales.

The complete results for all five architectures can be found on our github repository in the jupyter notebook *11\_multiplex\_enrichment*.

### Communities and transcriptome-wide association studies of disease

Leveraging the communities may have important methodological implications on the search for disease-associated genes. We asked whether incorporation of the communities would improve our ability to detect significant gene-level (TWAS/PrediXcan) associations (see “Methods” section). We chose to perform a TWAS of CRP in 361,194 UK Biobank subjects, as CRP is a biomarker of chronic low-grade inflammation, with elevated CRP levels associated with

a broad array of complex diseases, including cardiovascular disease, Alzheimer’s disease, and schizophrenia<sup>25</sup>. Notably, we found a significantly greater enrichment for associations with CRP (defined as adjusted  $p < 0.05$ ) among the set of genes that belong to a community than among the complement set of genes. In particular, the genes in communities showed a greater departure from the null (expected) distribution than the complement set of genes (Fig. 7a). This observation suggests that the use of the communities can substantially improve the signal-to-noise ratio in TWAS even in the case where the dataset is already highly-powered to detect causal associations. The estimated true positive rate  $\hat{\pi}_1$  (see “Methods” section) for association with the trait for the set of genes in communities was 0.45 while the estimate for the complement set was 0.37. Our top association with CRP among the community-located genes was *OASL* ( $p = 1.56 \times 10^{-55}$ ; see Fig. 7d for the chromosomal positions of the top associations), which has been previously implicated as a CRP and cardiovascular disease associated gene<sup>26</sup>. A similar performance gain in TWAS was observed for other traits (e.g., hemoglobin concentration and white blood cell [leukocyte] count for whole blood gene expression; see Supplementary Material and Supplementary Figs. 8 and 9), further demonstrating generalizability.

### Variational autoencoder model of communities and phenotypic consequences

Methodologically, the communities may also enable discovery of biologically-meaningful features in high-dimensional molecular data. We implemented a variational autoencoder (VAE) model (see “Methods” section) of the communities in 11,060 samples across 33 different cancer types in the TCGA data, customizing the *Tybat*



approach<sup>8</sup>. One benefit of a VAE is that it offers a probabilistic model, allowing us to do inference on the latent variable  $z$ , i.e.,  $P(z|X)$ . The marginal log-likelihood,  $\log P(X)$ , is generally intractable, which poses a challenge to this inference; however, Eq. (14) provides a variational lower bound on this marginal. In the VAE model, we assumed that the approximating posterior distribution  $Q(z|X)$  is multivariate Gaussian, whose mean  $\mu(X)$  and diagonal covariance matrix  $\sigma^2(X)I$  are learned by a neural network, so as to leverage the "reparametrisation trick"<sup>27</sup>. The second stage, which is also implemented as a neural network, generates a reconstructed representation of  $X$  from the stochastic  $z$ . In our implementation, we randomly split the input into a training set (80%) and a test set (20%). The encoded layer is compressed into a vector of size 100 consisting of a mean and variance. We assumed a learning rate of 0.0005, 50 epochs, and batch size of 50. We used Rectified Linear Unit (ReLU) activation for the encoder and sigmoid activation for the decoder. The optimization algorithm *Adam* was applied in the training to minimize the VAE loss, i.e., the negative of the Evidence Lower Bound Objective (ELBO) given by Eq. (14). The VAE loss as a function of epoch number (showing the training performance) and the reconstruction accuracy for the communities are largely equivalent to the corresponding results for the full transcriptome (Supplementary Fig. 10) in the TCGA data.

Notably, the latent representations learned by the VAE model from the communities encode biologically-meaningful features. The model learned to stratify metastatic tumors from primary tumors and normal solid tissues (Fig. 7b) (Mann–Whitney  $U$ -test  $p < 2.2 \times 10^{-16}$ ). This held robustly after adjusting for race, sex, age at diagnosis, or stage (logistic regression  $p < 2.2 \times 10^{-16}$  for each) or disease ( $p = 8.9 \times 10^{-5}$ ). Cancer progression may be characterized by oncogenic dedifferentiation (i.e., steady loss of differentiated phenotype) and acquisition of stemness (i.e., self-renewal and generation of differentiated progeny). The VAE model learned a stemness index that was significantly associated (Spearman  $\rho \approx -0.27$  in  $\log_2$  transformed space,  $p < 2.2 \times 10^{-16}$ ), across the spectrum of tumor types, with a recently developed DNA-methylation based index (Fig. 7c)<sup>28</sup> (see "Methods" section). Again, adjustment for race, sex, age at diagnosis, stage, or disease did not affect the result (linear regression  $p < 2.2 \times 10^{-16}$ ).

## DISCUSSION

We developed an inter-tissue multiplex framework for the analysis of the human transcriptome. Given the complexity of pathophysiological processes underlying complex diseases, intra-tissue and inter-tissue transcriptome analysis should enable a more complete mechanistic understanding. For these phenotypes, studying the interaction among tissues may provide greater insights into disease biology than an intra-tissue approach. Communities in co-expression networks were here shown to be enriched for some known pathways, encoding current understandings of biological processes; however, we identified other communities that are likely to contain novel or previously inaccessible functional information. Methodologically, use of the communities in TWAS and a neural network demonstrated substantial gain in the identification of disease-associated genes and discovery of biologically-meaningful information.

UMAP embeddings of the transcriptome of the entire set of communities (representing only 18% of all genes) fully revealed the tissue clusters. Low-dimensional representation of the subset of communities that are in the multiplex networks did not recover the tissue clusters, but uncovers other clustering patterns, suggesting a hierarchy of clusters at increasingly finer scales. We developed an approach to quantify the conservation of, and uncertainty in, the UMAP global structure and estimated the sampling distribution of the local structure (e.g., distances among tissue clusters), with broad relevance to other applications such as

cell population identification in single-cell transcriptome studies. New gene expression data can be embedded into our models, facilitating integrative analyses of the large volume of transcriptome data that are increasingly available. Notably, in external TCGA data, UMAP representations of the genes that belong to a GTEx-derived community induced clustering by cancer type, demonstrating the cross-study relevance of our approach. We provide a publicly available resource of co-expression networks, communities, multiplex architectures, enriched pathways, and code to stimulate research into network-based studies of the transcriptome.

Using the global multiplexity index, we investigated the tissue-sharedness of identified communities. In fact, communities that are shared across multiple tissues may suggest the presence of a tissue-to-tissue mechanism, that controls the activity of member genes across the layers in the network. Such regulatory mechanisms have been relatively understudied in comparison with intra-tissue controls.

We identified tissue-dependent communities that are enriched for human diseases. For example, we found a 15-member community in the "Brain Hippocampus" that is enriched for diseases associated with glycosaminoglycan metabolism. These genetic disorders are due to mutations in the biosynthetic enzymes for glycosaminoglycans, such as glycosyltransferases and sulfotransferases. Sulfated glycosaminoglycans include the chondroitin sulfate and dermatan sulfate chains that are covalently bound to the core proteins of proteoglycans, which are present in the extracellular matrices and at cell surfaces. Mutations affecting the biosynthesis of these chains may lead to genetic diseases that are characterized by craniofacial dysmorphism and developmental delay<sup>29</sup>. The community structure we identified proposes a cooperative role for these genes, and the fact that they span multiple chromosomes suggests the presence of coordinated transcriptional regulation.

Some of the communities are shared across multiple tissues; their dysregulation may thus lead to pleiotropic effects and contribute to known and novel comorbidities. We modeled these communities as belonging to layers of multiplex networks. For example, we identified a 17-member community in the "Brain Tissues" multiplex network (i.e., spanning across all brain regions sampled here), consisting primarily of ribosomal proteins. This community was enriched for proteins involved in translation (adjusted  $p = 2.53 \times 10^{-35}$ ), with significant overrepresentation for the viral mRNA translation pathway (adjusted  $p = 1.06 \times 10^{-38}$ ), NMD (adjusted  $p = 1.01 \times 10^{-37}$ ), and other Reactome pathways. Viral mRNAs in the cytoplasm can be translated by the host cell ribosomal translational apparatus, and indeed viruses have evolved strategies to recruit the host translation initiation factors necessary for the translation initiation by host cell mRNAs<sup>30</sup>. NMD, a surveillance pathway that targets mRNAs with aberrant features for degradation, may interfere with the hijacking of the host translational machinery<sup>31</sup>. In the brain, NMD, as a post-transcriptional mechanism, affects neural development, neural stem cell differentiation decisions, and synaptic plasticity; thus, defects in the pathway can cause aberrant neuronal activation and neurodevelopmental disorders<sup>22</sup>. Detecting this co-expression network of ribosomal proteins therefore provides a sanity check to our approach, but the identified community structure and the presence of this in the multiplex may suggest a highly coordinated regulatory mechanism across the tissues.

Leveraging the communities in TWAS of CRP in 361,194 subjects resulted in substantial performance gain in the discovery of trait-associated genes. Future methodological work that integrates community detection and additional omics data may further optimize the performance gain. A variational autoencoder implementation applied to the genes in the communities identified disease-relevant latent subspaces. Notably, this model learned a latent representation that significantly distinguishes

metastatic from primary tumors and another related to stem cell-associated molecular features, across the 11,060 samples in TCGA data. Prediction of metastasis and the biology of the stemness phenotype can be further investigated through the genes and their communities identified here, but more definitive conclusions will require extensive biological and clinical validation studies. Nevertheless, methodologically, a deep learning framework that explicitly exploits the topological structures in co-expression networks holds promise for uncovering critical biological insights into disease mechanisms, for example via integration of perturbations of the community topology or of the pleiotropic impact of communities shared across certain layers or tissues).

In summary, we performed network analysis on the most comprehensive human transcriptome dataset available to gain insights into how structures in co-expression networks may contribute to biological pathways and mediate disease processes. The rich resource we generated and the network approach we developed may prove useful to other omics datasets, facilitating studies of inter-tissue and intra-tissue regulatory mechanisms, with important implications for our mechanistic understanding of human disease.

## METHODS

### GTEx dataset

The GTEx V8 dataset<sup>11,15</sup> is a genomic resource consisting of 948 donors and 17,382 RNA-Seq samples from 52 tissues and two cell lines. The resource provides a catalog of genetic effects on the transcriptome and a broad survey of individual- and tissue-specific gene expression. Of the 54 tissues and cell lines, 49 include samples with at least 70 subjects, forming the basis of the analysis of genetic regulatory effects<sup>11</sup>. In this study, we leveraged the 49 tissues because of their sample size (see Supplementary Table 2) and our interest in shared transcriptional regulatory programs for co-expressed genes.

### Data preprocessing

We restricted our analyses to protein-encoding genes based on the GENCODE Release 26 (GRCh38) annotation. Although the GTEx dataset had annotated genes with Ensembl IDs, we removed duplicates (using GENE IDs) and unmapped genes from downstream analyses. After this preprocessing step, the resulting dataset is characterized by the following count statistics:

- Unique genes across all tissues: 18,364.
- Genes present in only 1 tissue: 412.
- Genes present in all 49 tissues: 12,557.

### Accounting for unmodelled factors

In order to correct for batch effects and other unwanted variation in the gene expression data, we used the *sva* R package (v3.34.0), which is specifically targeted for identifying surrogate variables in high-dimensional data sets<sup>32</sup>. For each tissue gene expression matrix, the number of components (latent factors) was estimated using a permutation procedure, as described by Buja and Eyuboglu<sup>33</sup>.

Subsequently, using the function *sva\_network*, residuals were generated after regressing out the surrogate variables. The residual values, rather than the original gene expression values, were used in the downstream analyses. For convenience, we refer to the residual values as the “gene expression data”, since they represent the expression levels that have been corrected for (unwanted) confounders.

### Tissue-dependent correlation and adjacency matrices

For each tissue, a correlation matrix  $C = [z_{ij}]$  was created by calculating the Pearson correlation coefficient  $r_{ij}$  for every pair  $(i, j)$  of genes. Fisher z-transformation was then applied:

$$z_{ij} = 0.5 \times \ln \left( \frac{1 + r_{ij}}{1 - r_{ij}} \right), \quad (1)$$

where  $\ln$  is the natural logarithm function.

For each correlation matrix, we retained only the strongest correlations (i.e., transformed  $z_{ij}$  less than  $-0.8$  and greater than  $0.8$ ) to generate a co-expression network. An adjacency matrix  $A = [A_{ij}]$  was defined, for each tissue, such that  $A_{ij}$  is equal to  $z_{ij}$  if gene  $i$  and gene  $j$  are co-expressed (retained), and zero otherwise. We assumed undirected networks without self-loops, which implies  $A_{ij} = A_{ji}$  and  $A_{ii} = 0$ .

### Community detection

We sought to detect groups of genes in each tissue, with the aim of finding communities whose internal connections are denser than the connections with the rest of the co-expression network. We applied the Louvain community detection method<sup>34</sup> in each tissue to generate a comprehensive atlas of communities. An asymmetric treatment for the negative correlations was used, thus inducing negatively correlated genes to belong to different communities<sup>35</sup>. The algorithm identifies communities by maximizing the modularity index<sup>36</sup>,  $Q$ , as the algorithm progresses:

$$Q^* = \frac{1}{v^+} \sum_{ij} (w_{ij}^+ - e_{ij}^+) \delta_{M_i M_j} - \frac{1}{v^+ + v^-} \sum_{ij} (w_{ij}^- - e_{ij}^-) \delta_{M_i M_j}. \quad (2)$$

Here, a positive connection between nodes  $i$  and  $j$  is denoted as  $w_{ij}^+$  and has a value between 0 and 1; likewise, a negative connection is represented  $w_{ij}^-$  and can also have a value between 0 and 1.  $e_{ij}^\pm$  is the chance-expected within-module connection weight and calculated, for each positive/negative correspondent, as  $\frac{s_i^+ s_j^+}{v^+}$ , where  $s_i^\pm$  is the sum of positive or negative connection weights of node  $i$ .  $v^\pm$  is the sum of all positive or negative edges, and  $\delta_{M_i M_j} = 1$  when nodes  $i$  and  $j$  are in the same module or zero otherwise. In particular, the Louvain method initially assigns each node to its own community and iteratively evaluates the gain in modularity, if one node is moved from one formed community to another of its neighborhood. We leveraged the Brain Connectivity Toolbox Python package v0.5.0 (available on github: [aestrivex/bctpy](https://github.com/aestrivex/bctpy)). The resolution parameter  $\gamma$  was set to its default value, 1.

### UMAP embeddings of community-defined gene expression

To produce a lower dimensional representation of the original dataset, we applied Uniform Manifold Approximation and Projection (UMAP)<sup>12</sup>, a manifold learning technique. Our goal was to generate a map that reveals embedded structures and test whether biologically relevant clusters can be recovered from the gene expression data. Towards this end, we analysed both the full master matrix  $\mathbf{M}$  of scaled gene expression (in the range  $[0, 1]$ ), consisting of all genes (i.e., 18,364), and a submatrix consisting of only those genes that belong to a community in at least one tissue (i.e., 3259). (Similarly to all of the results in the rest of the paper, we considered only Louvain communities with at least four genes.)

We chose UMAP because of the substantial improvement in running time on our data (compared to t-SNE, with its known computational and memory complexity that is quadratic in the sample size<sup>37</sup>), and UMAP’s theoretical grounding in manifold theory<sup>12</sup>. UMAP can also capture non-linear effects in gene expression, and this was another reason why we chose it, over more traditional dimensionality reduction techniques such as principal component analysis. Additional implementation details can be found in Supplementary Fig. 2.

### Persistence of the UMAP global structure

We quantified the conservation of, and variability in, the UMAP structure, including the relation among biologically-meaningful clusters, e.g., tissues. We characterized such a structure using the matrix  $[d(i, j)]$  of pairwise distances for clusters  $i$  and  $j$  in  $\{1, 2, \dots, L\}$ . For the actual (original) gene expression data, we define  $\mathbf{V}_{(0)} = [v_{ij,0}]$  as the resulting matrix of pairwise distances. Note  $\mathbf{V}_{(0)}$  is a symmetric matrix with zeros along the diagonal. We sought to:

- estimate the sampling distribution of  $d(i, j)$ , and calculate its standard error and a confidence interval,
- correlate the matrix  $\mathbf{V}_{(0)}$  and the resulting matrix  $W$  from a perturbation of the original structure.

We approached the quantification problem through a (non-parametric) bootstrapping procedure. From the master matrix  $\mathbf{M}$  of gene expression, we generated a total of  $B$  bootstrapped manifolds, each of equal size (here, each such sample was randomly drawn from 80% of the data points, i.e., rows, in  $\mathbf{M}$ ). For the  $k$ -th sample, we constructed the matrix  $\mathbf{V}_{(k)} = [d(i, j)_{(k)}]$  of pairwise distances derived from the UMAP embeddings for

tissues  $i$  and  $j$ . Here we used the "induced metric"  $d: \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  from the embedding  $\phi: M_g \hookrightarrow \mathbb{R}^m$  of the Riemannian manifold  $M_g$  into Euclidean space, but our treatment here generalizes to the intrinsic metric  $g: M_g \times M_g \rightarrow \mathbb{R}$ , with  $g(\phi^{-1}(i), \phi^{-1}(j))$ , of the original manifold. The set  $\{\widehat{d(i,j)}_{(k)}\}_{k=1}^B$  allows us to calculate the mean and variance of the UMAP-derived estimator for  $d(i, j)$ :

$$\overline{d(i,j)} = \frac{\sum_{k=1}^B \widehat{d(i,j)}_{(k)}}{B}, \quad (3)$$

$$\widehat{\sigma}_{d(i,j)}^2 = \frac{\sum_{k=1}^B \widehat{d(i,j)}_{(k)}^2}{B-1} - \left( \frac{\sum_{k=1}^B \widehat{d(i,j)}_{(k)}}{B-1} \right)^2. \quad (4)$$

This approach provides a maximum likelihood estimate, i.e.,  $\widehat{\sigma}_{d(i,j)}$ , of the standard error. We used a heatmap to visualize  $\overline{d(i,j)}$  for each tissue pair  $(i, j)$ . An alternative could have been to use a normalized "metric" (which is more robust to the scale from the embedding  $\phi$ ):

$$d^*(i,j)_{(k)} = \frac{d(i,j)_{(k)}}{\text{median}_{s,t \in 1, \dots, L} d(s,t)_{(k)}}, \quad (5)$$

but we found this normalization to be unnecessary in the GTEx data.

For two tissues  $i_0$  and  $i_1$ , we define a "clustering conservation coefficient" to quantify the preservation of the clustering of tissues  $i_0$  and  $i_1$  relative to all tissues  $\{j\}$ :

$$C_{(i_0, i_1)} = \text{corr}(\overline{d(i_0, j)}, \overline{d(i_1, j)}), \quad (6)$$

where  $\text{corr}$  is the correlation operator. The correlation is calculated for a pair of UMAP-derived distance estimates across all tissues  $\{j\}$ . In particular, this statistic allows us to formally test the null hypothesis of no conservation of global structure for a given pair of tissues under the null hypothesis,

$$\sqrt{\frac{L-3}{1.06}} \text{arctanh}(C_{(i_0, i_1)}) \sim N(0, 1). \quad (7)$$

This coefficient can be extended to a larger set of tissues,  $i_0, \dots, i_j$  (e.g., the 13 brain regions), using the first order statistic:

$$C_{i_0, \dots, i_j} = \min_{s, t \in 1, \dots, L} C_{(i_s, i_t)}. \quad (8)$$

Furthermore, we calculated the relationship between the original  $\mathbf{V}_{(0)}$  and "perturbed"  $\mathbf{V}_{(k)}$  for each sample  $k$ :

$$r_k = \text{corr}(\mathbf{V}_{(0)}, \mathbf{V}_{(k)}), \quad (9)$$

and the resulting empirical distribution of the correlation values  $r_k$ . We note that UMAP has a stochastic element since it utilizes stochastic approximate nearest neighbor search and stochastic gradient descent for optimization; however, the  $r_k$  derived from a different run  $\mathbf{V}_{(k)}$  (rather than from bootstrapping) quantifies the stability of the global structure in the presence of stochasticity. Collectively, our approach provides a way to perform statistical inference on the UMAP embedded structures.

### Prediction power of communities for tissues

We investigated the extent to which each community's gene expression profile was predictive of each of the tissues. The master matrix  $\mathbf{M}$ , representing the entire dataset under analysis, has 15,201 rows representing each RNA-Seq sample from each tissue collected from all subjects, and 18,364 columns representing the total number of genes available. If a value was non-existent (which may be due to the gene's expression being tissue-specific), we assumed a zero value, conveying no expression in that tissue.

For each community, the expression values of the member genes were selected from  $\mathbf{M}$ . With this sliced table, 49 binary classifications were performed using Support Vector Machine (SVM), wherein for each classification, we predicted each tissue. Essentially, the sliced table, which comprises the training data, for a  $k$ -member community can be viewed as a collection of vectors  $\{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$ , where  $\vec{x}_i \in \mathbb{R}^k$  is the gene expression profile of the  $k$  genes for the  $i$ -th sample and  $y_i \in \{1, 0\}$  indicates membership in the tissue to be predicted. The goal of the classification is to separate the tissue to be predicted from the other tissues via the largest margin hyperplane, which can be generically written as  $\vec{w} \cdot \vec{x} + b = 0$ , where  $\vec{w}$  is normal to the hyperplane. SVM was used with a linear kernel and weights were adjusted to be inversely proportional to class frequencies in the input data (this corresponds to setting the *class\_weight*

parameter in *scikit-learn* to "balanced"). To avoid overfitting, each classification was performed using a stratified 3-fold cross-validation procedure, in which the  $F_1$  score metric was used to report the prediction power across the three folds. We decided to use the  $F_1$  score instead of other metrics, given that each binary classification was highly unbalanced, i.e., a given tissue is the positive outcome and all the other 48 tissues are the negative outcome. (Louvain communities with less than four genes were filtered out from this analysis.)

$$F_1 = \frac{2}{(\text{Precision})^{-1} + (\text{Recall})^{-1}} \quad (10)$$

### Enrichment analysis

To evaluate the degree to which a community corresponds to well-known biological pathways, we performed enrichment analyses using the Reactome 2016 as reference. We used the *gseapy* python package to make calls on the *Enrichr* web AP<sup>38</sup>. As per the *Enrichr* official documentation, the  $p$ -value is computed using Fisher's exact test (hypergeometric test). We considered significant those pathways with a Benjamini-Hochberg-adjusted  $p$ -value below 0.05. Louvain communities with less than four genes were considered "not enriched".

### Multilayer analysis

In order to investigate the tissue-shared profiles of gene communities, as well as the relationships between gene expression traits across tissues, we proceeded to model our system as a multilayer network<sup>39</sup>. Formally, a multilayer network is defined as a pair  $\mathbf{A} = (\mathbf{G}, \mathbf{D})$ , where  $\mathbf{G} := \{G_1, \dots, G_L\}$  is a set of graphs and  $\mathbf{D}$  consists of a set of interlayer connections existing between the graphs and connecting the different layers. Each graph  $G_l \in \mathbf{G}$  is a "network layer" with its own associated adjacency matrix  $A_l$ . Thus,  $\mathbf{G}$  can be specified by the vector of adjacency matrices of the  $L$  layers:  $\mathbf{A} := (A_1, \dots, A_L)$ . Multilayer networks allow us to represent complex relationships which would otherwise be impossible to describe using single-layer graphs separately considered.

A special case of multilayer networks is a multiplex network, which we used to model the GTEx transcriptome data. In this case, all layers are composed of the same set of nodes but may exhibit very different topologies. The degree of node  $i$  is the vector  $d^{[i]} = (d_1^{[i]}, \dots, d_L^{[i]})$ , and  $d_1^{[i]}$  may vary across the layers. Interlayer connections are established between corresponding nodes across different layers. Layers represent different tissues, nodes represent genes, and edges between two nodes are weighted according to the correlation weights. In the GTEx data, the correlation matrices, previously described, define an adjacency matrix  $A_l$  for each layer  $l$  of the multiplex network.

Using the communities of co-expressed genes for each tissue, we then computed the so-called *global multiplexity index*<sup>40</sup> to investigate the relationships of communities across different layers. This index quantifies how many times two nodes (genes) are clustered in the same communities across different layers. If, for example, gene  $i$  and gene  $j$  are clustered together in the layer of tissue  $T_1$  and of tissue  $T_2$ , then the global multiplexity index is two. In the matrix  $\text{gmi}(i, j)$  of global multiplexity indices for a multiplex architecture, each element represents the number of times that two given genes,  $i$  and  $j$ , are clustered in the same community. More formally, if  $L$  is the number of layers,  $N$  the number of nodes for each layer, and  $c_l^g$  the community membership of gene  $i$  at graph  $g$ , then the global multiplexity index  $\text{gmi}(i, j)$  for gene  $i$  and gene  $j$ , with  $i$  and  $j \in \{0, \dots, N\}$  is defined as follows:

$$\text{gmi}(i, j) = \sum_{g=1}^L \delta(c_l^g, c_l^j), \quad (11)$$

where  $\delta(c_l^g, c_l^j)$  represents the Kronecker delta function. The value of  $\text{gmi}(i, j)$  therefore increases by 1 if the two nodes are found to be part of the same community in a layer. If two genes share a high value of global multiplexity index, this may indicate a greater level of connectivity and suggest greater functional similarity, as they appear multiple times in the same community across different layers. We define the individual probabilities:

$$\pi_k = p(\text{gmi} = k \mid \theta), \quad (12)$$

where  $\theta$  represents all of the parameters of the model.  $\pi_k$  is the probability that two genes are clustered in the same communities across  $k$  layers, where  $k \leq L$  and  $L$  is the number of layers in the network. We estimated the probability distribution of  $\text{gmi}(i, j)$  in the GTEx data.

We tested whether the UMAP embeddings of the transcriptome profiles of the communities in a multiplex architecture—a subset of all communities previously interrogated—could also recover biologically-meaningful clusters. This analysis allowed us to estimate the topology of the high-dimensional transcriptome data and test whether additional clusters could be uncovered at increasingly finer scales.

### Application to transcriptome-wide association studies

To evaluate the relevance of the communities for genomic studies of human disease, we performed TWAS/PrediXcan analysis of C-reactive protein (CRP)<sup>5,7</sup>. We chose CRP for its clinical significance as a biomarker for a wide range of complex diseases, including cardiovascular disease, type 2 diabetes mellitus, Alzheimer's disease, and age-related macular degeneration<sup>41</sup>. Briefly, TWAS/PrediXcan estimates the "genetically-determined component of gene expression" in genome-wide association study (GWAS) subjects and infers the gene's association with the phenotype. The inference can be done using GWAS summary statistics<sup>6,18</sup>. We hypothesized that prioritization of the communities in TWAS/PrediXcan would improve the signal-to-noise ratio for detecting gene-level associations. We applied whole blood (local genetic variation based) gene expression prediction models<sup>5</sup> trained in GTEx v8 data<sup>42</sup> to a GWAS of CRP in 361,194 samples (of white-British ancestry) in the UK Biobank<sup>43</sup> (nealelab.is). We compared the associations derived from the set of genes that belong to a community and from the complement set of genes using a conditional Quantile–Quantile (Q–Q) plot (i.e., conditional on community membership status) of empirical quantiles of nominal negative  $\log_{10}(p)$  values. The conditional Q–Q plot for each set of genes can be framed in terms of the false discovery rate (FDR)<sup>44</sup>; at a  $p$ -value threshold, the Bayes FDR is given by:

$$\text{FDR}(p) = \pi_0 F_0(p) / F(p), \quad (13)$$

where  $\pi_0$  is the proportion of null genes,  $F_0$  is the null cumulative distribution function (cdf), and  $F$  is the cdf for both null and non-null genes. Under the null hypothesis,  $F_0$  is the cdf of the standard uniform distribution (i.e.,  $F_0(p) = p$  for  $p \in [0, 1]$ ) while  $F$  can be estimated by the empirical distribution. We estimated the true positive rate for each of the two non-overlapping sets of genes defined by community membership using  $\pi_1 := 1 - \pi_0$ , as previously described<sup>18</sup>.

### Variational autoencoder (VAE) model of communities

We implemented a VAE<sup>27</sup>, a deep learning methodology, to learn biologically-meaningful latent representations of the transcriptome<sup>8</sup> or a subset, e.g., the genes in the communities. VAE is a two-phase generative model, and we implemented it to capture major sources of variation with non-linear effects. In the encoding phase, dimensionality reduction is performed on the input; the decoding phase performs reconstruction of the original input from a latent and stochastic representation. The following equation serves as the basis for the VAE:

$$\log P(X) - \mathcal{D}[Q(z|X)||P(z|X)] = E_{z \sim Q}[\log P(X|z)] - \mathcal{D}[Q(z|X)||P(z)], \quad (14)$$

where  $P(X)$  is a probability density defined for a data point  $X$  in the input expression data (e.g., the set of genes that belong to the communities),  $z$  is a vector of latent variables with a prior probability density function  $P(z)$ , and  $\mathcal{D}$  is the Kullback–Leibler (KL) divergence between  $P(z|X)$  and some function  $Q(z)$ . The first part of the right hand side of the equation, i.e., the expected negative log-likelihood, gives the reconstruction loss while the second is the KL divergence between the learned latent distribution and the prior distribution.  $Q(z|X)$  is the probabilistic encoder which compresses the data  $X$  into the latent variable  $z$ , whereas the generative model  $P(X|z)$  is the probabilistic decoder which reconstructs the latent representation into the original data.

Using a VAE model of the communities, we tested the extent to which they could help to identify expression changes associated with disease and discover biologically-meaningful features<sup>8,45</sup>. We built on *Tybolt*, which uses a *Keras* implementation. We analysed TCGA Pan-Cancer data, consisting of batch-effects-normalized mRNA data (in units of  $\log_2(\text{norm\_value} + 1)$ ) in 11,060 samples across 33 cancer types (see "Data and code availability" section). Expression values were mapped to [0,1] for each gene using the maximum and minimum values.

We evaluated the performance of the VAE model of the GTEx-derived communities in TCGA data in two ways. (1) The process of metastasis remains poorly understood. Classification of the tumors into primary or metastatic origin enabled us to test whether the VAE model could

successfully discriminate the sample type. (2) The acquisition of a stem cell-like tumor trait, i.e., stemness, in cancer suggests gene expression programs that may contribute to progression and treatment resistance. To determine whether the VAE model successfully learned stemness<sup>28</sup>, we used a DNA methylation-based "Stemness Score" from the Pan-Cancer Stemness working group, specifically the "DNAs" signature, which combines (a) an epigenetically-regulated DNA methylation-based signature, (b) a differentially-methylated probes-based signature, and (c) an enhancer elements/DNA methylation-based signature. For these analyses, we also adjusted for race, sex, age at diagnosis, stage, or disease as a potential confounding effect.

### Statistical tests

All statistical tests were two-sided.

### DATA AVAILABILITY

The protected data for the GTEx project (for example, genotype and RNA-sequence data) are available via access request to dbGaP accession number phs000424.v8.p2. Processed GTEx data (for example, gene expression and eQTLs) are available on the GTEx portal: <https://gtexportal.org>. TCGA gene expression data and DNA methylation based stemness scores can be downloaded from the University of California, Santa Cruz (UCSC) TCGA Pan-Cancer Atlas hub on Xena: <https://pancanatlas.xenahubs.net>. The UK Biobank C-Reactive Protein GWAS summary results on which TWAS/PrediXcan was applied are available for download: <https://doi.org/10.5281/zenodo.4681322>. We also deployed the results on Track Hub of the UCSC Genome Browser: [https://genome-euro.ucsc.edu/cgi-bin/hgTracks?hgid=261097667\\_sHbF4mmwke2A6qjRyKmNiVAot6m](https://genome-euro.ucsc.edu/cgi-bin/hgTracks?hgid=261097667_sHbF4mmwke2A6qjRyKmNiVAot6m).

### CODE AVAILABILITY

Code for reproducibility and results are publicly available on Github, with additional instructions for implementation: <https://github.com/tjiagoM/gtex-transcriptome-modelling>.

Received: 12 November 2020; Accepted: 28 April 2021;

Published online: 27 May 2021

### REFERENCES

- Fortunato, S. Community detection in graphs. *Phys. Rep.* **486**, 75–174 (2010).
- Saha, A. et al. Co-expression networks reveal the tissue-specific regulation of transcription and splicing. *Genome Res.* **27**, 1843–1858 (2017).
- Gerring, Z. F., Gamazon, E. R. & and, E. M. D. A gene co-expression network-based analysis of multiple brain tissues reveals novel genes and molecular pathways underlying major depression. *PLoS Genet.* **15**, e1008245 (2019).
- Gamazon, E. R. et al. Using an atlas of gene regulation across 44 human tissues to inform complex disease-and trait-associated variation. *Nat. Genet.* **50**, 956–967 (2018).
- Gamazon, E. R. et al. A gene-based association method for mapping traits using reference transcriptome data. *Nat. Genet.* **47**, 1091 (2015).
- Gusev, A. et al. Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.* **48**, 245 (2016).
- Zhou, D. et al. A unified framework for joint-tissue transcriptome-wide association and mendelian randomization analysis. *Nat. Genet.* **52**, 1239–1246 (2020).
- Way, G. P. & Greene, C. S. Extracting a biologically relevant latent space from cancer transcriptomes with variational autoencoders. *Bioinformatics* **2018**, 80–91 (2017).
- Eraslan, G., Simon, L. M., Mircea, M., Mueller, N. S. & Theis, F. J. Single-cell rna-seq denoising using a deep count autoencoder. *Nat. Commun.* **10**, 1–14 (2019).
- Leek, J. T. & Storey, J. D. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet.* **3**, e161 (2007).
- The GTEx Consortium. The GTEx consortium atlas of genetic regulatory effects across human tissues. *Science* **369**, 1318–1330 (2020).
- McInnes, L., Healy, J., Saul, N. & Grossberger, L. UMAP: Uniform manifold approximation and projection. *J. Open Sourc. Softw.* **3**, 861 (2018).
- Diaz-Papkovich, A., Anderson-Trocme, L. & Gravel, S. Umap reveals cryptic population structure and phenotype heterogeneity in large genomic cohorts. *PLoS Genet.* **15**, e1008432 (2019).
- Rijnberk, A. In *Clinical Endocrinology of Dogs and Cats* 11–34 (Springer, 1996).
- Consortium, G. et al. Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).

16. Margolis, R., Margolis, R., Chang, L. & Preti, C. Glycosaminoglycans of brain during development. *Biochemistry* **14**, 85–88 (1975).
17. Huynh, M. B. et al. Glycosaminoglycans from alzheimer's disease hippocampus have altered capacities to bind and regulate growth factors activities and to bind tau. *PLoS ONE* **14**, e0209573 (2019).
18. Gamazon, E. R., Zwinderman, A. H., Cox, N. J., Denys, D. & Derks, E. M. Multi-tissue transcriptome analyses identify genetic mechanisms underlying neuropsychiatric traits. *Nat. Genet.* **51**, 933–940 (2019).
19. Rogers, G. et al. From gut dysbiosis to altered brain function and mental illness: mechanisms and pathways. *Mol. Psychiatry* **21**, 738–748 (2016).
20. Li, G.-c et al. Identification of novel biomarker and therapeutic target candidates for acute intracerebral hemorrhage by quantitative plasma proteomics. *Clin. Proteom.* **14**, 14 (2017).
21. Siems, S. B. et al. Proteome profile of peripheral myelin in healthy mice and in a neuropathy model. *Elife* **9**, e51406 (2020).
22. Jaffrey, S. R. & Wilkinson, M. F. Nonsense-mediated rna decay in the brain: emerging modulator of neural development and disease. *Nat. Rev. Neurosci.* **19**, 715–728 (2018).
23. Rao, M. & Gershon, M. D. The bowel and beyond: the enteric nervous system in neurological disorders. *Nat. Rev. Gastroenterol. Hepatol.* **13**, 517 (2016).
24. Majewski, M., Kozłowska, A., Thoene, M., Lepiarczyk, E. & Grzegorzewski, W. Overview of the role of vitamins and minerals on the kynurenine pathway in health and disease. *J. Physiol. Pharmacol.* **67**, 3–19 (2016).
25. Ligthart, S. et al. Genome analyses of >200,000 individuals identify 58 loci for chronic inflammation and highlight pathways that link inflammation and complex disorders. *Am. J. Hum. Genet.* **103**, 691–706 (2018).
26. Middelberg, R. P. et al. Genetic variants in lpl, oasl and tomm40/apoe-c1-c2-c4 genes are associated with multiple cardiovascular-related traits. *BMC Med. Genet.* **12**, 123 (2011).
27. Kingma, D. P. & Welling, M. Auto-encoding variational bayes. Preprint at <https://arxiv.org/abs/1312.6114> (2013).
28. Malta, T. M. et al. Machine learning identifies stemness features associated with oncogenic dedifferentiation. *Cell* **173**, 338–354 (2018).
29. Mizumoto, S., Ikegawa, S. & Sugahara, K. Human genetic disorders caused by mutations in genes encoding biosynthetic enzymes for sulfated glycosaminoglycans. *J. Biol. Chem.* **288**, 10953–10961 (2013).
30. Bushell, M. & Sarnow, P. Hijacking the translation apparatus by rna viruses. *J. Cell Biol.* **158**, 395–399 (2002).
31. Rigby, R. E. & Rehwinkel, J. Rna degradation in antiviral immunity and autoimmunity. *Trends Immunol.* **36**, 179–188 (2015).
32. Parsana, P. et al. Addressing confounding artifacts in reconstruction of gene co-expression networks. *Genome Biol.* **20**, 1–6 (2019).
33. Buja, A. & Eyuboglu, N. Remarks on parallel analysis. *Multivar. Behav. Res.* **27**, 509–540 (1992).
34. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech.* **2008**, P10008 (2008).
35. Rubinov, M. & Sporns, O. Weight-conserving characterization of complex functional brain networks. *NeuroImage* **56**, 2068–2079 (2011).
36. Newman, M. E. & Girvan, M. Finding and evaluating community structure in networks. *Phys. Rev. E* **69**, 026113 (2004).
37. Maaten, Lvd & Hinton, G. Visualizing data using t-sne. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).
38. Kuleshov, M. V. et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **44**, W90–W97 (2016).
39. Kivelä, M. et al. Multilayer networks. *Journal of Complex Netw.* **2**, 203–271 (2014).
40. Hristova, D., Rutherford, A., Anson, J., Luengo-Oroz, M. & Mascolo, C. The international postal network and other global flows as proxies for national wellbeing. *PLoS ONE* **11**, e0155976 (2016).
41. Luan, Y.-y & Yao, Y.-m The clinical significance and potential role of c-reactive protein in chronic inflammatory and neurodegenerative diseases. *Front. Immunol.* **9**, 1302 (2018).
42. Barbeira, A. N. et al. Exploiting the GTEx resources to decipher the mechanisms at GWAS loci. *Genome Biol.* **22**, 1–24 (2021).
43. Bycroft, C. et al. The uk biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
44. Andreassen, O. A. et al. Improved detection of common variants associated with schizophrenia and bipolar disorder using pleiotropy-informed conditional false discovery rate. *PLoS Genet.* **9**, e1003455 (2013).
45. Maj, C. et al. Integration of machine learning methods to dissect genetically imputed transcriptomic profiles in alzheimer's disease. *Front. Genet.* **10**, 726 (2019).

## ACKNOWLEDGEMENTS

E.R.G. is grateful to Clare Hall, University of Cambridge for the Fellowship support. We thank The Genotype-Tissue Expression (GTEx) project for the use of v8 data and The Cancer Genome Atlas (TCGA) project for the use of gene expression data from the cancer samples. This research has been conducted using the UK Biobank Resource. This research is supported by the National Institutes of Health (NIH) Genomic Innovator Award R35HG010718, NIH/NHGRI R01HG011138, NIH/NIGMS R01GM140287, and NIH/NHLBI R01HL133559. T.A. is funded by the W. D. Armstrong Trust Fund, University of Cambridge, UK.

## AUTHOR CONTRIBUTIONS

T.A. and G.M.D. conducted the analysis. E.R.G., T.A. and G.M.D. wrote the manuscript. E.R.G. and P.L. designed and supervised the study. All authors contributed to the editing of the manuscript. T.A. and G.M.D. contributed equally to this work.

## COMPETING INTERESTS

E.R.G. receives an honorarium from the journal Circulation Research of the American Heart Association, as a member of the Editorial Board. He performed consulting on pharmacogenetic analysis with the City of Hope/Beckman Research Institute. The other authors declare no competing interests.

## ADDITIONAL INFORMATION

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41540-021-00186-6>.

**Correspondence** and requests for materials should be addressed to P.L. or E.R.G.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021