

# ORANGE: A Machine Learning Approach for Modeling Tissue-Specific Aging from Transcriptomic Data

Wasif Jalal<sup>1</sup>, Mubasshira Musarrat<sup>1</sup>, Md. Abul Hassan Samee<sup>2</sup>, and M. Sohel Rahman<sup>1</sup>

<sup>1</sup>Department of Computer Science and Engineering, Bangladesh University of Engineering and Technology

<sup>2</sup>Department of Integrative Physiology, Baylor College of Medicine

## Abstract

Despite aging being a fundamental biological process which profoundly influences health and disease, the interplay between tissue-specific aging and mortality remains underexplored. This study applies machine learning on GTEx transcriptomic data to model tissue-specific biological ages across 12 different types of tissues and introduces an *age-gap* metric to quantify deviations from the chronological age. Our best models achieve an average RMSE of 6.44 years and an average  $R^2$  of 0.64. Age-gap statistics reveal significant tissue-specific aging patterns, identifying extreme agers and correlations between extreme aging and mortality. About 20% of subjects are found to exhibit extreme aging in one tissue, while 1% show multi-organ aging. These findings greatly emphasize the role of transcriptomics in aging research and its implications for health and longevity.

## 1 Introduction

Aging is closely tied to the onset of numerous diseases and the overall health span of humans. Recent advancements in molecular biology and machine learning have enabled researchers to delve deeper into the intricacies of biological aging, uncovering organ-specific patterns that diverge from chronological aging. A notable study by Oh et al. (2023) [29] utilized plasma proteomic data obtained via SomaScan [11] [40] assays to predict organ-specific biological ages and introduced the concept of *age-gaps*. These age-gaps, defined as the difference between an individual's chronological age and the biological age predicted from proteomic data, have been shown to correlate with health and disease states, offering a new lens through which to study the interplay of aging and pathology. Inspired by this approach, we sought to explore the application of similar methodologies to transcriptomic data, specifically gene expression values measured as transcripts per million (TPM) from tissue samples in the GTEx dataset, incorporating subject sex as an additional variable. Inspired by the findings of Oh et al. [29] regarding the prediction of organ-specific age from the plasma proteome, we aimed to explore the potential of predicting tissue-specific age from the human

transcriptome using machine learning models.

Despite the several epigenetic and proteomic studies in illuminating aging processes, the potential of transcriptomic data to predict organ-specific biological age remains underexplored. The GTEx dataset, a comprehensive resource of gene expression profiles across 54 human tissue types, provides an opportunity to investigate this avenue. However, unlike plasma proteomics, transcriptomic data is inherently tissue-specific and requires a different approach to account for the unique expression patterns of each tissue. Additionally, challenges, such as limited sample sizes per tissue type, variability in donor characteristics, and the lack of precise chronological age data in the public GTEx dataset complicate the modeling process. Addressing these challenges necessitates innovative preprocessing, feature selection, and predictive modeling techniques.

A recent study by Johnson and Krishnan (2023) [16] has associated the transcriptome with age and sex, utilizing transcriptomic data from RNA-seq samples to predict the age group of subjects within each sex group. However, it did not analyze the data on a tissue-specific basis, leaving unexplored the unique expression patterns that may exist across different tissue types. Furthermore, a multimodal approach to age estimation was recently applied to ovary and lung tissues from the adult GTEx dataset [35]. This approach combined transcriptomic data with additional modalities to improve prediction accuracy. For lung tissues, we improved upon the RMSE and  $R^2$  scores of their elastic net-based transcriptomic predictors. Remarkably, the RMSE and  $R^2$  scores of our models are competitive even when compared to their multimodal ensembles making our simpler models more attractive and acceptable in this context. Previously, Ren and Kuan (2020) [34] instituted methods to predict tissue-specific and cross-tissue ages using the transcriptomic data from GTEx. While their methodology provided a foundation for transcriptomics-based age prediction, our models are simpler, and in most tissues, demonstrate improved performance in terms of RMSE and  $R^2$  scores. Furthermore, we extend the scope of transcriptomic age prediction by incorporating the concept of *age-gaps* and related statistics inspired by Oh et al. (2023) [29], along with their associations with mortality.

This integration of statistical analyses of age-gaps introduces a novel perspective on tissue-specific aging and its broader implications for health and disease.

Building upon prior research, we developed machine learning models to predict organ-specific biological ages using transcriptomic data from 12 selected tissues, focusing on tissues with sufficient samples and known correlations with mortality. Our methodology included robust feature selection methods, such as identifying age-correlated genes and analyzing differentially expressed genes with age, alongside advanced modeling techniques like bootstrap-aggregated Partial Least Squares (PLS) regression. By introducing a transcriptomics-based *age-gap* metric, we aimed to assess deviations in predicted biological ages and their associations with health outcomes, mortality, and potential tissue-specific aging regulators. The results of our study highlight the viability of using transcriptomic data for organ-specific age prediction. Our models demonstrated competitive performance, achieving an average RMSE of approximately six years and uncovering distinct patterns of extreme aging in specific tissues. We observed notable statistics in our analysis of age-gaps. Nearly 20% of the population in our study exhibited strongly accelerated aging in one organ, while a smaller subset, approximately 1%, was identified as multi-organ agers. These accelerated aging patterns were associated with mortality, conferring around 3 times the risk of death from intermediate or prolonged illness.

In this study, we provide a novel application of the *age-gap* concept to transcriptomic data, expanding its utility beyond proteomics. Moreover, our study underscores the importance of tissue-specific analysis in aging research, using the GTEx dataset to develop robust predictive models. Lastly, we identify transcriptomic markers associated with tissue aging, laying the groundwork for future studies exploring the molecular mechanisms of aging and disease. Through this work, we bridge the gap between transcriptomics and aging research, offering a new perspective on the dynamic interplay between biological age, health, and disease.

## 2 Methods

### 2.1 Dataset

We use the **Adult Genotype-Tissue Expression (GTEx)** v10 [24] dataset, which provides gene expression data measured in transcripts per million (TPM) across 54 tissue types from 948 adult subjects [7]. The dataset contains 19,788 tissue samples, encompassing various tissues, such as, brain, liver, lungs, colon, pancreas, arteries, and heart. This extensive dataset, collected from post-mortem tissue samples, includes gene expression data along with metadata [8] on age (binned into 10-year ranges), sex, circumstance of death characterized by the Hardy Scale (see Supplementary Section 1.1 for details on this scale) [28], and other clinical details, making it a valuable resource for understand-

ing tissue-specific gene expression and genetic variations across a diverse human population.

### 2.2 Preprocessing and Initial Analysis

Initial analysis involved generating t-SNE plots of the tissue-specific gene TPMs, resulting in distinct clustering of samples corresponding to individual tissues (Figure S2).

#### 2.2.1 Tissue Selection for Predictive Modeling

Since the Adult GTEx does not contain samples of every tissue type from each test subject, and the aging of all organs does not have a strong correlation with mortality, we had to limit our study to a set of tissues that had a significant number of samples and are known to be correlated with mortality [4] [27]. The selected tissue types were liver, aorta, coronary artery, brain cortex, brain cerebellum, heart atrial appendage, subcutaneous adipose, lung, sun-exposed skin, tibial nerve, sigmoid colon, and pancreas. We plotted t-SNE plots for the samples for individual tissue types and observed no distinct clustering based on age range. However, in certain cases, clustering was evident when categorized by the Hardy Scale ratings (Figure S3).

#### 2.2.2 Optimal fixed-point interpolation for age-ranges

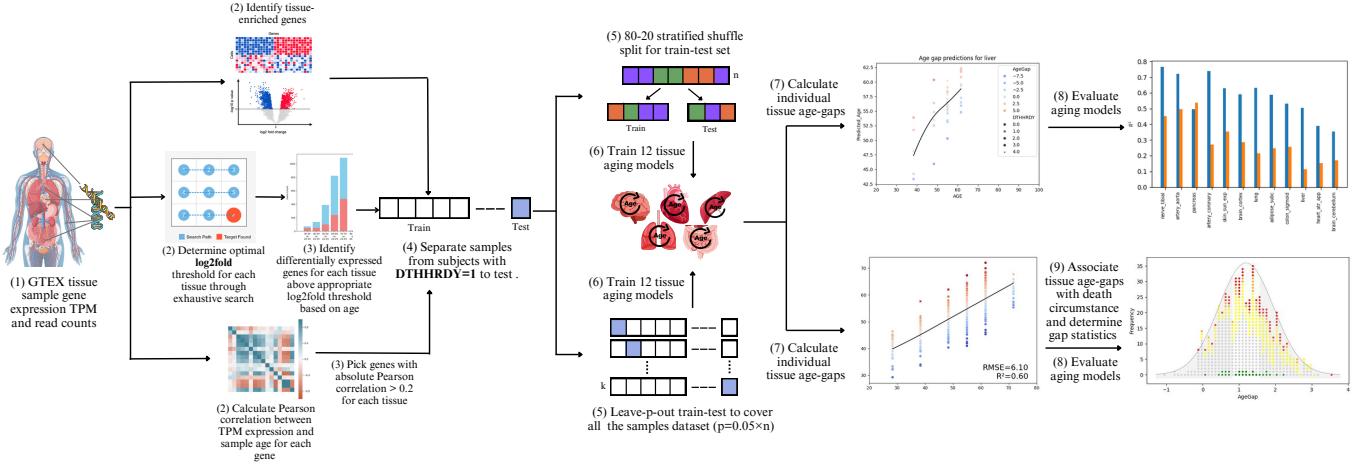
To compensate for the absence of precise ages of the subjects in the dataset, we opted to find optimal points within each age range using an exhaustive search approach, rather than just taking the midpoint of each range. Each age range was divided into three points using the formula:

$$Age = Mid-Point + [(Ternary\ Digit - 1) \times \frac{Range}{3}] \quad (1)$$

Iterating through  $3^6$  permutations on six age ranges for each tissue, we found the most commonly-occurring optimal permutation of ternary digits to be 222110. As expected from a Gaussian distribution (Figure S1), the optimal point for each age range deviated from the range mid-values towards the population mean.

### 2.3 Predicting Age from Gene Expression

To develop tissue-specific aging models, we grouped the gene expression data by tissue type, and then performed all training and testing on tissue-specific data only. We initially partitioned each tissue's dataset into training and testing subsets, ensuring that all samples with a Hardy Scale rating of 1 (Supplementary Section 1.1) were allocated to the testing sets, since the age at death for subjects who died of unnatural causes such as accident or suicide may not be reflective of their tissue health. We



**Figure 1: Study Methodology.** The study utilized GTEx v10 data to predict tissue-specific biological aging from gene expression profiles. Key tissues were selected based on sample size and mortality correlations, and machine learning models were trained on age-correlated features for each tissue. An *age-gap* metric was introduced to quantify deviations in predicted tissue age from chronological age, revealing extreme agers in specific tissues. The analysis highlighted patterns of accelerated aging in a small subset of subjects, providing insights into tissue-specific aging dynamics.

assigned the remaining subjects to the training and testing subsets randomly to get an overall train-test ratio of 80:20.

### 2.3.1 Feature selection by identifying age-correlated tissue-specific genes

We selected features for tissue aging models applying three different methods. In the first method, we conducted a correlation analysis on the columns of expression values for each tissue, and selected only the columns exhibiting a Pearson correlation greater than 0.2, thereby retaining features that demonstrate a significant linear relationship. In the second method, we chose genes that were differentially-expressed with age, using the package PyDESeq2 [26], which is a Python implementation of the Bioconductor package DESeq2 [25]. We fixed the optimal  $\log_2(fold-change)$  threshold for each tissue using an exhaustive search approach, and selected the genes that showed the highest fold-changes across age groups. In the third method, we followed an approach almost identical to the one used by Oh et al. [29] to identify organ-enriched plasma proteins. A gene is organ-enriched if it is expressed at least four times higher in a single organ compared to any other organ [43]. We define tissue-enriched genes as genes that are expressed at least four times higher in a single tissue than in any other tissue. The  $\log_2(fold-change)$  of normalized read count of each gene in each tissue was determined using the PyDESeq2 [26] package, and each gene was assigned as a feature to the tissue in which its fold-change was the greatest.

### 2.3.2 Feature transformation

In order to stabilize variance and make the gene expression data more Gaussian-like, the Yeo-Johnson power transformation (Eqn. 2) was applied to the expression

value columns.

$$y(\lambda) = \begin{cases} \frac{((y+1)^\lambda - 1)}{\lambda} & \text{if } \lambda \neq 0, y \geq 0 \\ \ln(y+1) & \text{if } \lambda = 0, y \geq 0 \\ \frac{-((-y+1)^{2-\lambda} - 1)}{2-\lambda} & \text{if } \lambda \neq 2, y < 0 \\ -\ln(-y+1) & \text{if } \lambda = 2, y < 0 \end{cases} \quad (2)$$

### 2.3.3 Learning algorithms to predict sample age

Utilizing the selected gene expression columns for each organ, along with sex (encoded as 0/1) and the optimal fixed point within each age range from the age column as target variables, we trained linear predictors. The results were averaged through bootstrap aggregation to enhance model robustness and improve predictive accuracy.

For modeling tissue-specific age, we opted to employ 20x bootstrapped predictors of several types. Initially, following the method used by Oh et al. [29], we tried LASSO (Least Absolute Shrinkage and Selection Operator) regression, which prevents overfitting using L1 regularization. The LASSO regression minimizes the following objective:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^n (y_i - X_i^T \beta)^2 + \alpha \sum_{j=1}^p |\beta_j|, \quad (3)$$

where  $X$  is the feature matrix,  $y$  is the target vector,  $\beta$  represents the regression coefficients,  $n$  is the number of observations,  $p$  is the number of features, and  $\alpha$  is the regularization parameter controlling the amount of shrinkage.

We also considered elastic net regularization, which combines L1 and L2 regularization at a ratio of 0.5 alongside L1 regularization. Elastic net regression minimizes the following objective:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^n (y_i - X_i^T \beta)^2 + \alpha \left( \frac{1-\rho}{2} \sum_{j=1}^p \beta_j^2 + \rho \sum_{j=1}^p |\beta_j| \right), \quad (4)$$

where  $\rho$  controls the balance between L1 and L2 regularization, and  $\alpha$  remains the overall regularization strength.

We then applied Partial Least Squares (PLS) regression [46], which projects the features and targets onto a shared low-dimensional space to model their relationship effectively. This method is particularly advantageous when predictors are highly collinear or when the number of predictors exceeds the number of observations. PLS regression creates new predictor variables, called components, as linear combinations of the original predictor variables, selected to maximize the covariance with the target variable. These components are then used to predict the response variable. The PLS model assumes a latent variable approach. The observed variables  $X$  and  $y$  are described as linear combinations of unobserved latent variables plus random noise as follows:

$$X = TP^T + E, \quad y = Tq + f, \quad (5)$$

where  $T$  is the score matrix,  $P$  is the loading matrix for  $X$ ,  $q$  is the loading vector for  $y$ ,  $E$  and  $f$  are residuals, and  $T$  is shared between  $X$  and  $y$ , ensuring that the components represent a shared structure.

In parallel to the above linear predictors, we also experimented with Support Vector Regression (SVR) [9] with linear kernels. SVR aims to find a hyperplane that predicts  $y$  with at most  $\epsilon$ -deviation while maximizing the margin. The objective is:

$$\min_{w,b,\xi,\xi^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*), \quad (6)$$

subject to

$$\begin{aligned} y_i - (w^T x_i + b) &\leq \epsilon + \xi_i, \\ (w^T x_i + b) - y_i &\leq \epsilon + \xi_i^*, \\ \xi_i, \xi_i^* &\geq 0. \end{aligned}$$

Here,  $C$  is the regularization parameter,  $\epsilon$  defines the margin of tolerance, and  $\xi, \xi^*$  are slack variables.

To explore the effectiveness of non-linear models, we tried using random forest regressors [3]. Random forests are ensemble methods that aggregate predictions from multiple decision trees to improve robustness and reduce overfitting. The prediction is given by:

$$\hat{y} = \frac{1}{n_{estimators}} \sum_{k=1}^{n_{estimators}} \hat{y}_k, \quad (7)$$

where  $\hat{y}_k$  is the prediction from the  $k$ -th decision tree and  $n_{estimators}$  is the number of trees in the ensemble.

### 2.3.4 Age-gap estimation

Following the prediction of tissue-specific ages, we calculated a regression for the predicted ages of tissue samples against the chronological ages of their respective subjects using a neural network containing two fully-connected layers ( $1 \times 32$  and  $32 \times 1$ ) with batch normalization and ReLU activation layers in between, thus yielding a  $\hat{y}$  value for each test sample. We used this approach instead of a simple linear regression to account for any possible nonlinearity in the underlying relationship between gene expression and age. Inspired by the concepts presented by Oh et al. [29], we defined a metric called the *age-gap*, by subtracting the  $\hat{y}$  value from the predicted age. This age-gap serves as an indicator of a tissue sample's aging relative to its counterparts within the same age group.

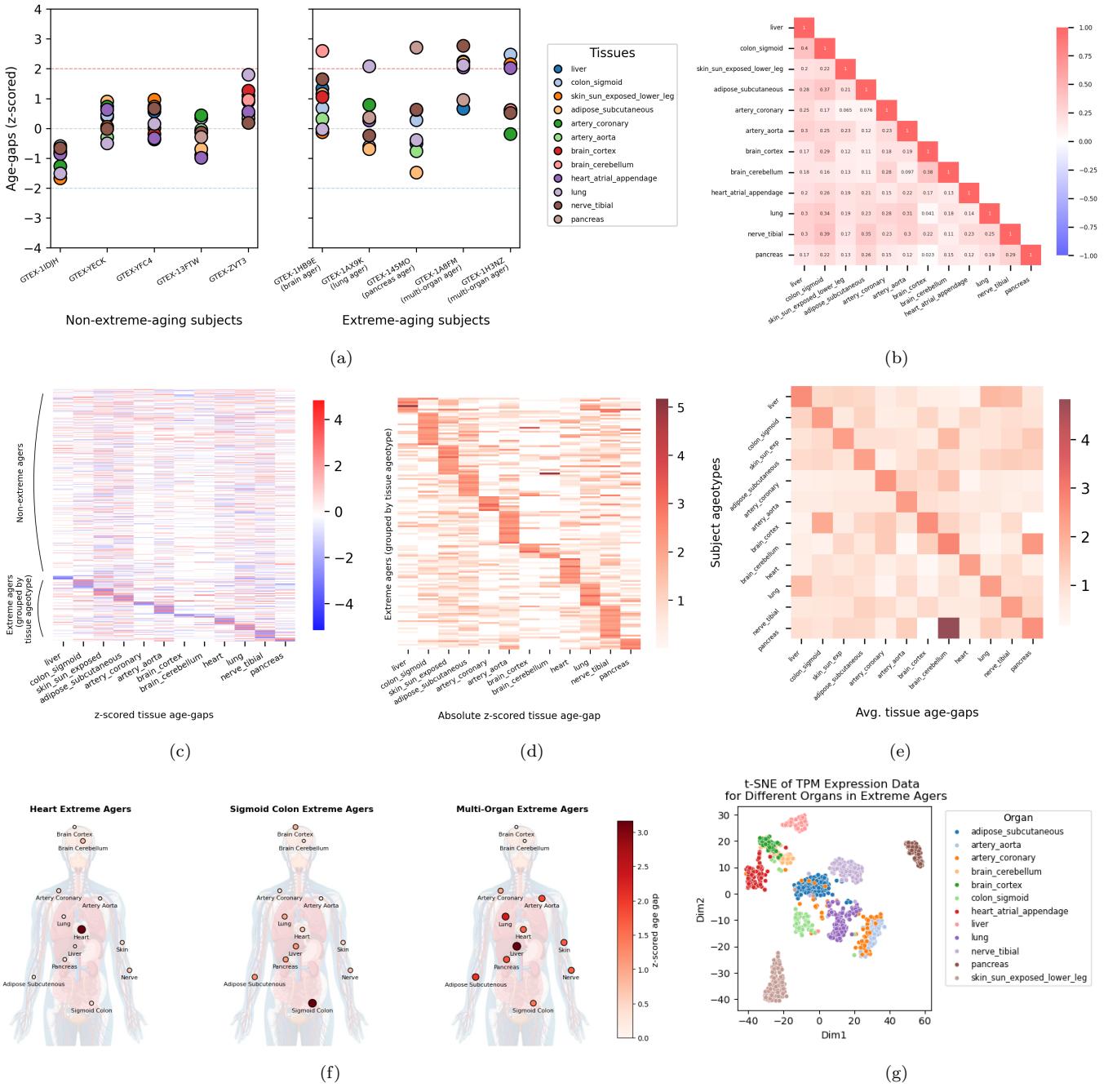
## 2.4 Result Interpretation and Downstream Analyses

### 2.4.1 Evaluation metrics

We evaluated the predictive performance of our models (i.e., predicted age vs. chronological age) using two metrics: the RMSE (Root mean-squared error) and the  $R^2$  (coefficient of determination). To choose the optimal model for downstream analyses, we also considered training speed as a metric.

### 2.4.2 Leave- $P$ -out training-testing for better downstream analyses on age-gaps

We aimed to generate a larger set of predicted age-gaps across all subjects, in order to examine the correlation between extreme aging and circumstance of death (characterized by Hardy scale ratings) through a conditional probability analysis. To compensate for the small subject count of the dataset, apart from the train-test splitting strategy presented in Section 2.3, we also adopted another strategy, where after allocating the samples with Hardy Scale ratings of 1 directly to the test set, we selected  $P$  of the shuffled samples from each tissue to perform a leave- $P$ -out train-test splitting in a sliding window mechanism, where  $P$  equals 5% of the total number of samples from that tissue. This strategy ensured that we covered each sample in the tissue's dataset. Each 5% window was used as the testing set, on which predictions were performed by training on the rest of the samples. This enabled us to perform a prediction on every single sample and ultimately have a larger set of predicted age-gaps, on which we could perform better downstream analyses. The samples with Hardy Scale ratings of 1, that we had initially separated for testing only, were appended to the last testing window of each tissue's dataset. For this leave- $P$ -out workflow and all downstream analyses, we used a training architecture that involved feature selection based on absolute Pearson correlation, followed by bootstrap-aggregated PLS regression without cross-validation. This approach provided the best balance between prediction accuracy and execution speed (see Sup-



**Figure 2: Distribution of extreme agers across the population.** (a): Extreme agers (at least one tissue's z-scored age-gap  $> 2$ ) and non-extreme agers. Five randomly-sampled subjects without extreme aging are shown; all appear to have varied tissue age-gaps that are slightly correlated and thus scattered around a mean. About 26% of all subjects were found to have extreme aging, and most of them (21% of all subjects) exhibited extreme aging in only one organ (three example subjects shown), while only 0.94% of the subjects showed multi-organ aging (two samples shown). (b): Heatmap of pairwise correlation between the age-gaps of each tissue type, 0.21 on average. (c): Heatmap of z-scored age-gaps of all subjects, partitioned into non-extreme agers and extreme agers grouped by tissue ageotype. (d): Heatmap of absolute z-scored age-gaps of extreme agers only, grouped by tissue ageotype (tissue with highest age-gap). Most agers show relatively accelerated aging in only one organ, by a notable margin. (e): Average age-gaps in each tissue (column-wise) of the subjects of each tissue ageotype (row-wise). Each ageotype has the highest average age-gap in the corresponding tissue by a notable margin. Note that the subjects with extreme-aging aorta tissues also have a higher average coronary artery age-gap, indicating that there may be a common group of cardiovascular agers among them. (f): Average tissue age-gaps of subjects with extreme aging in heart atrial appendage and sigmoid colon tissues shown as examples. Multi-organ agers show extreme aging across tissues, whereas most others show extreme aging in only one tissue along with low absolute age-gaps in other tissues. (g): t-SNE plot of tissue gene expression values (TPM) of extreme agers. Expression values of the same tissue cluster together, as do those of related tissues.

plementary Tables 2 and 3). We ran 20 cycles of this leave-P-out train-test exercise by varying a random seed each time that affected only the initial shuffling that we performed to ensure class homogeneity.

#### 2.4.3 Conditional probability analysis

For each tissue type, we made dot plots of the age-gaps calculated from our predictions through the leave- $P$ -out exercise and tried to fit them under normal distributions. At the two extremes of such distributions, we identified groups that we refer to as ***extreme negative agers*** and ***extreme positive agers***, defined as those falling below ( $\text{mean} - \frac{1}{2} \times \text{stddev}$ ) and above ( $\text{mean} + \frac{1}{2} \times \text{stddev}$ ), respectively.

$$P(\text{extreme}^+) = \frac{n(\text{gap} > \mu_{\text{gap}} + \frac{\sigma_{\text{gap}}}{2})}{n(\text{subjects})}$$

$$P(\text{extreme}^-) = \frac{n(\text{gap} < \mu_{\text{gap}} - \frac{\sigma_{\text{gap}}}{2})}{n(\text{subjects})}$$

Conversely, individuals situated near the center of the distribution, approximately within ( $\pm \frac{1}{5} \times \text{stddev}$ ) around the mean, were categorized as ***average agers***.

$$P(\text{average}) = \frac{n(\mu_{\text{gap}} + \frac{\sigma_{\text{gap}}}{5} > \text{gap} > \mu_{\text{gap}} - \frac{\sigma_{\text{gap}}}{5})}{n(\text{subjects})}$$

To explore the contributions of extreme age-gaps in disease and mortality, we considered the conditional probability ratios –

$$\frac{P(dthhrdy \in \{3, 4\} | \text{extreme}^+)}{P(dthhrdy \in \{3, 4\} | \text{extreme}^-)} \quad (8)$$

$$\frac{P(dthhrdy \in \{3, 4\} | \text{extreme}^+)}{P(dthhrdy \in \{3, 4\} | \text{average})} \quad (9)$$

On the other hand, to determine how age-gaps may relate to subjects who died unnatural deaths but may have been otherwise healthy ( $dthhrdy = 1$ ), we considered the following ratios –

$$\frac{P(dthhrdy = 1 | \text{extreme}^-)}{P(dthhrdy = 1 | \text{extreme}^+)} \quad (10)$$

$$\frac{P(dthhrdy = 1 | \text{extreme}^-)}{P(dthhrdy = 1 | \text{average})} \quad (11)$$

In the above equations,  $dthhrdy$  is the Hardy Scale rating (Supplementary Section 1.1) of the circumstance of subjects' deaths.

#### 2.4.4 Code, Environment and Availability

All the predictors were trained and tested using the Python package scikit-learn [32]. The hyperparameters  $\alpha$  (L1 regularization parameter for LASSO and elastic net),  $C$  (regularization parameter for SVR), and

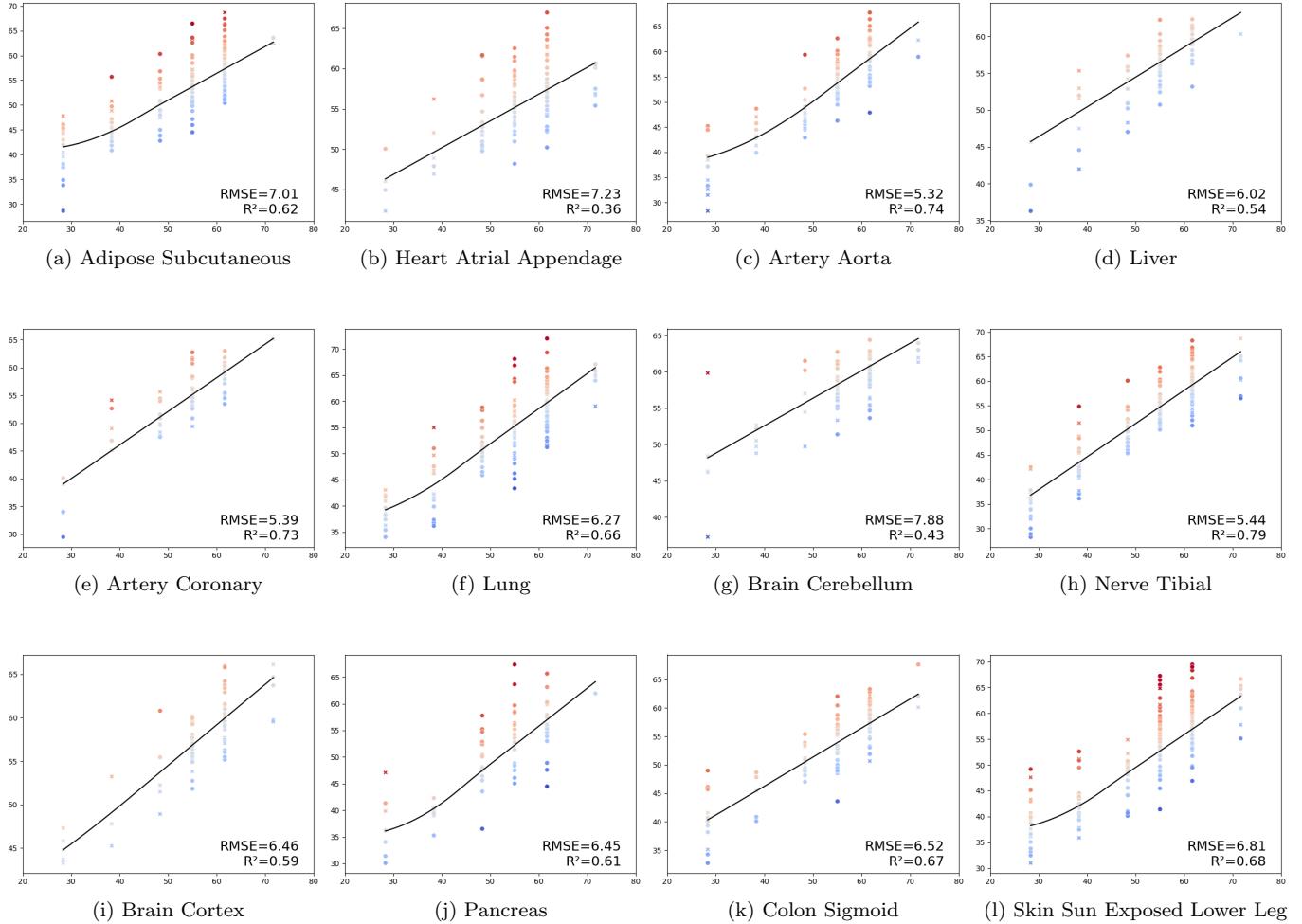
$n_{\text{estimators}}$  (number of trees in random forest regression) were all tuned with four-fold cross-validation using the GridSearchCV tool from scikit-learn. Fully-connected layers used in age-gap estimation were implemented using the PyTorch [30] library. We have made the code available through a GitHub repository (<https://github.com/wjalal/ORANGE/>).

## 3 Results

### 3.1 Model Performance

We observed varying model performance per tissue in terms of prediction accuracy and training speed in the various combinations of the methods we used for feature selection and training (Supplementary Tables 2 and 3). In general, the selection of characteristics by absolute Pearson correlation performed better than the selection by tissue-specific DEGs (differentially-expressed genes). Elastic net predictors had the best RMSE and  $R^2$  scores on average, with LASSO and PLS regressions closely following it. PLS regression without cross-validation performed the fastest while also predicting sample ages with competitive accuracy. The PLS approach greatly reduces the dimensionality of the dataset and performs much faster than the Ordinary Least Squares (OLS) regressors Lasso, Ridge, and Elastic Net which employ L1 and L2 regularization. PLS regression offered the best balance between (training) speed and accuracy. It also achieved competitive performance without the use of cross-validation in training, without which the OLS regressors could not perform as well.

When tested on common datasets, our PLS regression models also performed well compared to models developed in other studies (Supplementary Table 4). The RNAAgeCalc models developed by Ren et al. (2020) [34], which were trained on the GTEx v6 data, did not generalize well when we used them to predict the age of GTEx v10 samples that were not present in GTEx v6. However, our models, when trained on the portion of GTEx v10 data which overlaps with GTEx v6, achieved better average RMSE and  $R^2$  scores of 6.99 years and 0.49 respectively on test samples. Note that, the model was tested on GTEx v10 samples that were not present in GTEx v6. Our elastic net and PLS regression models for lung tissue aging also performed noticeably better than the gene expression-based aging models used by Ribeiro et al. (2024) [35] in their study on multimodal modeling of human aging. In fact, despite having a more constrained setup and having gene expression data as the sole modality, on a common testing set, our transcriptomic aging models exceeded the performance of all of their models except the epigenomic ones (methylation-based) and performed significantly better than their gene expression-based models (both elastic net and gradient boosted trees).



**Figure 3: Modeling tissue age with gene TPM.** Model performance across various tissue types, measured in  $R^2$  and RMSE, using Pearson correlation-based feature selection and elastic net predictors. Each plot shows the testing set samples of the corresponding tissue as dots whose positions along the plot’s x-axis and y-axis represent the corresponding sample’s chronological age and predicted age respectively. The line of fit in each plot is estimated by neural networks described in Section 2.3.4, while the vertical distance between each dot and the line of fit represents the corresponding sample’s age-gap.

### 3.2 Gene TPMs can model organ age

Our bootstrap-aggregated linear models could capture clear linear relationships between the TPM values of differentially-expressed genes measured from tissue samples and the subjects’ chronological ages (Figure 3). Our best models achieved an average (weighted by number of samples analyzed per tissue type) RMSE of 6.44, and  $R^2$  of 0.64 for predicted ages.

#### 3.2.1 Distribution of extreme agers across the population

We observed that most individuals had varying age-gaps across organs (Figure 2a). Analyzing the statistics of the predicted age-gaps among 531 subjects led us to some interesting findings. We observed low average correlation between the age-gaps of different tissues of a subject (Figure 2b), with an average pair-wise Pearson correlation of 0.21 across all tissue types. However, no negative correlation was observed, and in general, most of an individual’s

age-gaps (across tissue types) centered around an average (Figure 2a), due to a low correlation. About 26% of all subjects were found to have extreme aging (absolute value higher than two standard deviations) (Figure 2c). These gene expression values of the samples from these subjects clustered together by tissue, and the clusters of the individual tissues mostly remained separate from each other (Figure 2g). Notably, 21% of the subjects exhibited extreme aging in only one organ (Figure 2d), while only 0.94% of the subjects showed multi-organ aging (extreme aging in 3 or more tissues), suggesting that most cases of extreme aging are organ-specific (Figure 2e, 2f). With our novel approach, which is different from existing ones from the literature, we effectively reached conclusions that are in line with the current knowledge base [29]. In addition, we were able to identify some new relevant insights which call for further biological validation, thereby opening new avenues for research.

### 3.3 Age-gaps are associated with circumstance of death

Dot plots of each subject's tissue-specific age-gap or maximum observed age-gap across the selected 12 tissue types appear to fit normal distributions (Figure 4(a)-(l)). We could observe differences in the distribution of the five Hardy Scale death circumstances within the various regions under the normal curves (*extreme positive agers*, *extreme negative agers*, and *average agers*).

#### 3.3.1 Tissue age-gaps of vital organs correlate to cause of death

Across all 12 tissues that were analyzed, we found that *extreme positive agers* were significantly more likely to have died with  $dthhrdy = 3$  or  $4$  (intermediate death with a short terminal phase, or slow death after a long illness; both classified as ill subjects), than *extreme negative agers*, and moderately more likely compared to *average agers* (Figure 4m). We argue that the maximum age-gap that a subject has across the selected 12 tissue types, could be an indicator of their worst-aged organ, and possibly be associated with the subject's cause of death. From the conditional probability analysis of  $dthhrdy$  (Hardy Scale rating for circumstance of death) distribution by age-gap (Figure 4n, Supplementary Table 1), we could conclude that, *extreme positive agers* were 2.96 times as likely to have died with  $dthhrdy = 3$  or  $4$  (intermediate death with a short terminal phase, or slow death after a long illness; both classified as ill subjects), than *extreme negative agers*, and 1.24 times as likely as *average agers* (based on Eqns 8 and 9). It was also observed that *extreme negative agers* were 1.27 times as likely to have died with  $dthhrdy = 1$  (death due to accident, blunt force trauma, or suicide) as *extreme positive agers*, and 0.59 times as likely as *average agers* (based on Eqns 10 and 11), suggesting that subjects whose tissues had decelerated aging relative to their chronological age were more likely to have died due to unnatural causes rather than diseases.

### 3.4 Identifying important regulators for tissue aging

The highest-magnitude coefficients in our linear models for tissue age prediction offered insights into several genes that are known to specifically regulate aging in the corresponding tissues (Figure 5) as well as across multiple tissues. A total of 17 genes were identified as important age regulators across multiple tissues, while a total of 583 were observed as regulators of tissue-specific age.

#### 3.4.1 Tissue-independent or multi-tissue determinants of aging

Our aging models for each of the analyzed tissues, except for liver and brain, assigned high positive weights to the expression of PTCHD4, a gene which interacts

with pathways involved in DNA repair and cell cycle regulation, particularly in the context of the p53 pathway [6] [44]. PTCHD4-AS, a related long non-coding RNA, has shown potential in enhancing DNA repair and possibly activating checkpoints that could protect cells from accumulating damage- a factor that could influence aging by either supporting cell survival or promoting senescence to prevent malignancy [37].

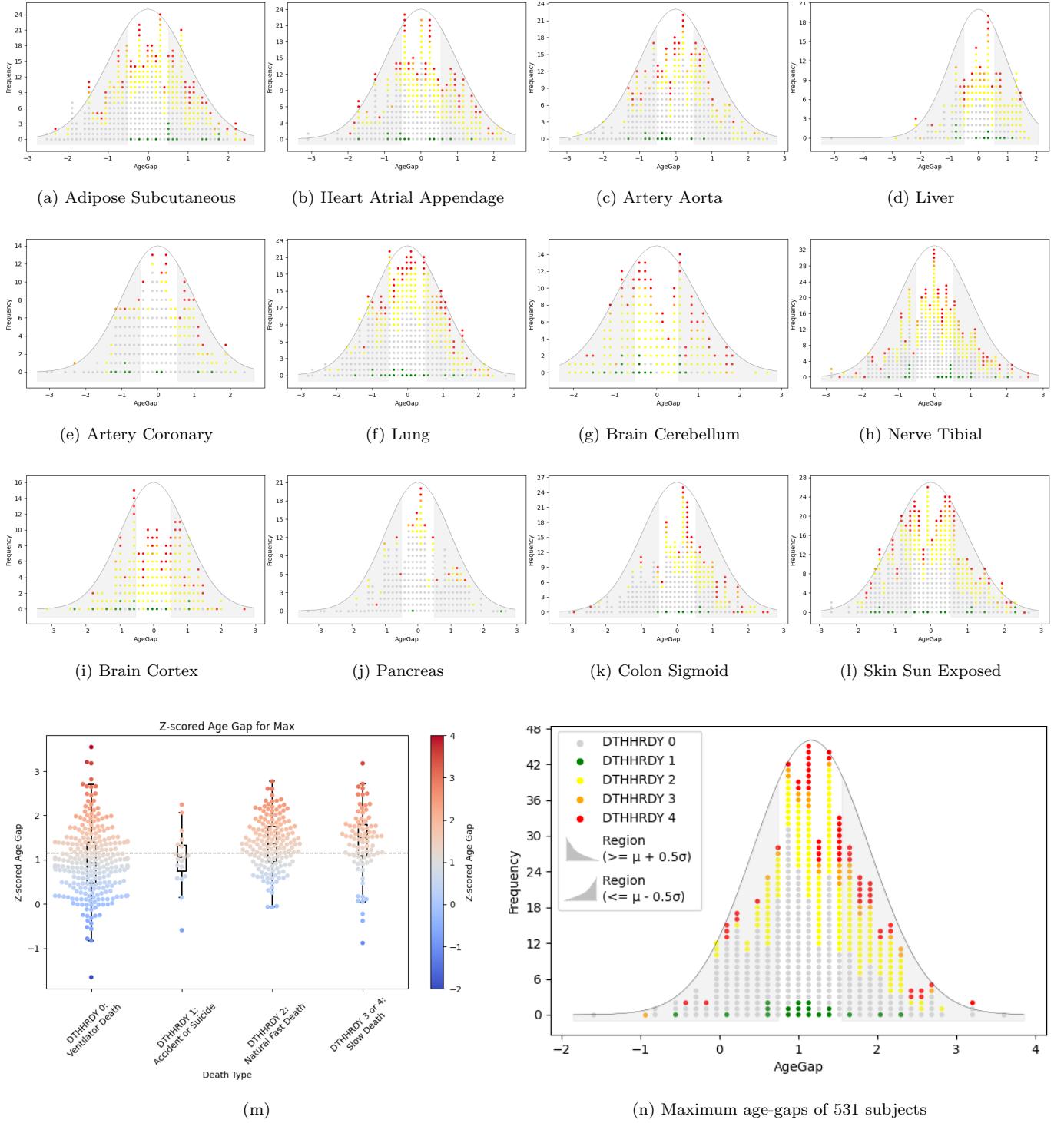
EDA2R, which has been described as the most remarkable and ubiquitous aging-related transcriptional hallmark [1] and a general marker of aging across tissues [34], was weighted significantly in our brain cortex, heart, adipose, lung, pancreas, and colon aging models. The expression of EDA2R in adipose, heart and lung was also independently found to be connected with aging [15] and cellular apoptosis, suggesting that it is an inducer of apoptosis [20]. LMO3, whose expression has been found to increase with age in adipose tissues, bears major weight in our artery, heart, adipose and skin tissue aging models. ZNF518B, which bears negative weight in several of our models, has been observed to be consistently downregulated with age in both humans and mice [39]. Interestingly, the pseudogene MTHFD2P1 is positively correlated with age in several of our tissue aging models, which prompts further biological studies to explore its role in aging, particularly because the related protein-coding gene MTHFD2 is called a potential oncogene [48] and its corresponding enzyme has been shown as an aging-associated factor in cancer [45].

MIR34AHG, the host gene for MicroRNA34-a, which is an important mediator of inflammaging [33], was also identified as a major contributor to age in several of our tissue aging models. Studies have shown that the expression of MicroRNA-34a increases with age [10] and it regulates cardiac aging and function [2], which has been reflected in the positive coefficients assigned to its expression in our models for heart atrial appendage and other tissues.

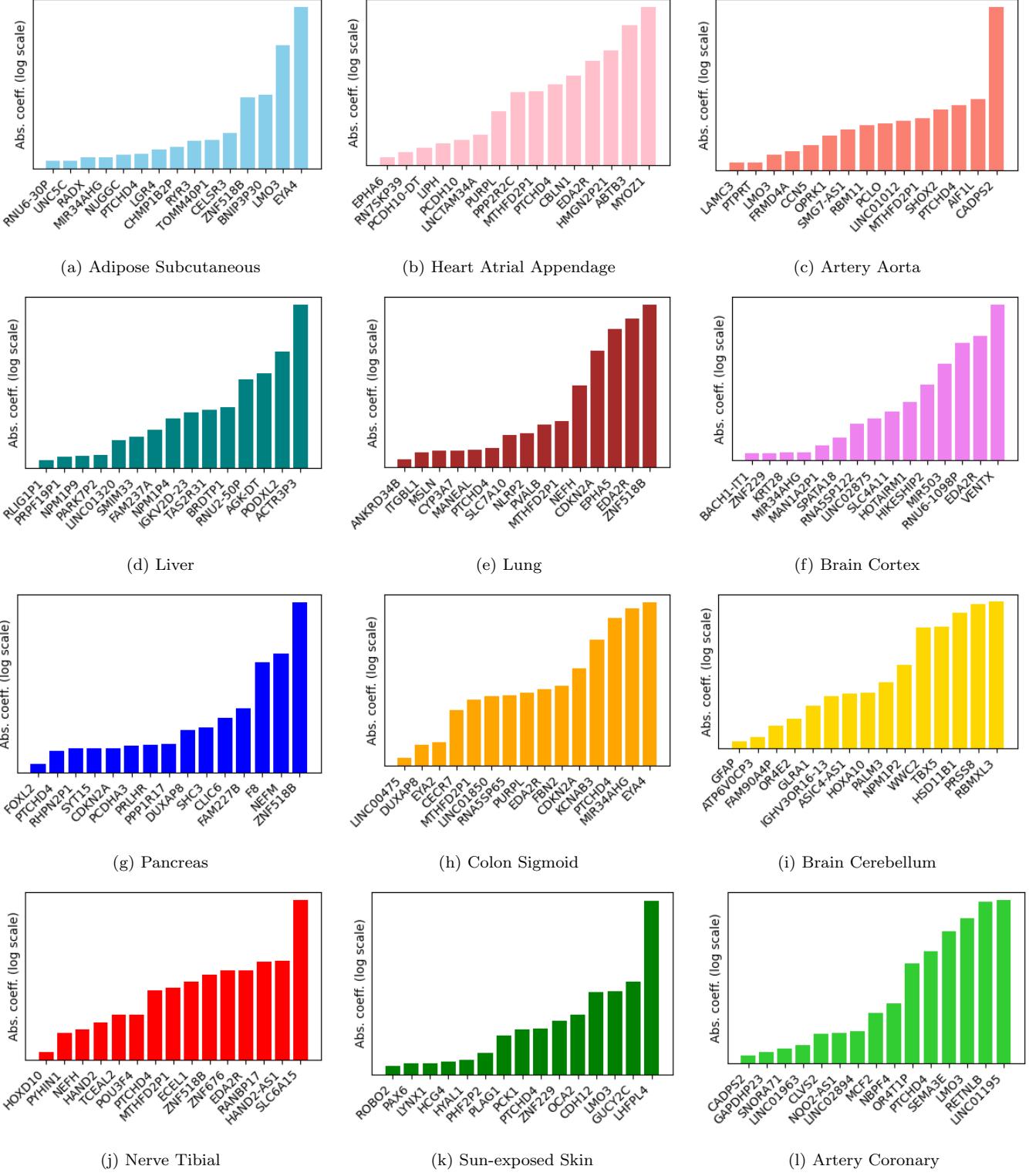
#### 3.4.2 Tissue-specific determinants of aging

One of the genes that our subcutaneous adipose tissue aging model identified as a key determinant of aging in the tissue was ZMAT3, whose upregulation has been associated with senescence in subcutaneous adipose [41], activation of the p53/p21 pathway, and the inhibition of adipogenesis. Another gene called OCA2 bears a negative weight in our skin aging model, thereby aligning with its established role in melanin production and skin pigmentation, reinforcing previous findings linking melanin levels to skin aging and cancer risk [12] [18].

The RNA gene, PURPL, which promotes tumorigenicity in colorectal cancer [23], was assigned a high coefficient by our sigmoid colon aging model. In addition, the model highlighted DUXAP8, a pseudogene



**Figure 4: Tissue age-gaps of vital organs correlate to cause of death.** (a)-(l): Distribution of age-gaps calculated by predicting ages of the samples of a tissue type. The age-gaps fit normal distributions. Regions of interest in the tails of the curve represent the conditions for **extreme negative agers** and **extreme positive agers**, defined as those falling below ( $\text{mean} - \frac{1}{2} \times \text{stddev}$ ) and above ( $\text{mean} + \frac{1}{2} \times \text{stddev}$ ), respectively. (m): Distribution of maximum age-gaps of 531 subjects with samples of more than 5 tissue types, compared across *dthhrdy* classes. Subjects who died of prolonged illness (*dthhrdy* 3 or 4) had their maximum (worst-case tissue) z-scored age-gaps distributed above the mean. (n): Distribution of maximum age-gaps of 531 subjects with samples of more than 5 tissue types. The legend in (n) applies to all the plots (a)-(l) and (n), to indicate the conditional probability of finding each *dthhrdy* type among the two types of extreme agers.



**Figure 5: Top 15 Absolute Coefficients for Various Tissues.** Plots show the significant features contributing to model predictions for each tissue.

that contributes to colorectal cancer progression by inducing epithelial-mesenchymal transition [13].

In the brain, our cerebellar aging model identified GFAP as a significant gene associated with aging. This gene encodes the glial fibrillary acidic protein, which is an intermediate filament protein found in astrocytes. Elevated GFAP levels in cerebrospinal fluid have been proposed as biomarkers for neurodegenerative diseases, such as Alzheimer’s disease [5], while GFAP has also been found to be elevated in the cerebral cortices of autistic subjects [21], suggesting that the encoding gene is a marker for neurodegeneration brain-aging. Additionally, our brain cortex aging model identified MIR34AHG as a key factor. MIR34A, a protein encoded by MIR34AHG, has been shown to modulate inflammatory molecules involved in post-stroke recovery and affect blood-brain barrier permeability [31]. Its potential role in Alzheimer’s disease, suggested by a study that finds it overexpressed in the cortex [38], further supports its relevance in brain aging.

Our lung aging models identified MSLN, the protein-coding gene for Mesothelin, as a regulator of aging, and its relevance is supported by several studies that link the expression of Mesothelin with human lung cancer [47] [14] [17]. EPHA5, another gene whose expression has been associated with the development of human lung adenocarcinoma, was also found to be an important regulator of lung aging by our model. On the other hand, NLRP2, which has been identified as an antioncogene which inhibits the proliferation of lung adenocarcinoma cells [22], was assigned a negative weight by our model.

In the pancreas, an important determinant of aging identified by our model is the tumor suppressor gene CDKN2A, which is a well-known catalyst in the development of pancreatic ductal adenocarcinoma [19]. Our pancreas aging model also assigns a strong weight to PRLHR which is the encoding gene for the prolactin-releasing hormone receptor. Its role in pancreas aging is supported by literature that establishes prolactin as a promoter of pancreatic cancer progression [42].

## 4 Discussions

In this study, we have developed a mechanism to model tissue-specific biological aging from the transcriptome of post-mortem tissue samples. Our models can estimate the biological age of tissue samples from the expression of identified genes measured as TPM (transcripts per million) through RNA-seq, model tissue-specific aging profiles, and identify transcriptomic markers of aging.

Although the GTEx is widely regarded as the largest dataset for tissue-specific human transcriptomic data, there were many limitations in our study owing to the restrictions of the dataset. Despite having a large

number of tissue samples, the number of individual subjects is around 1000, which is relatively small for analytics at the individual level. The open-access Adult GTEx only has limited data on donor phenotypes, which are sex, 10-year age bracket, and Hardy scale rating. Unbinned age, race, weight, smoking status, diabetes status, and other disease-related de-identified donor phenotypes are available in the protected-access dataset, which potentially holds future directions for our work. Despite using a 10-year age bracket, our age prediction models had RMSE scores around 6, which is comparable to or even better than models that used exact ages [34] [35]. It is also worth mentioning that the GTEx project collects samples only from tissue sites that it classifies as *non-diseased*; so the inclusion of diseased tissues could change the picture of modeling tissue-specific aging and disease risk. The project does however include samples from subjects with Hardy scale ratings defined as death after prolonged disease. Future studies on this topic should focus on tissue-specific disease prediction and modeling.

We had to restrict our study to a specific set of tissues, because not all the tissue-types in the dataset have sufficient samples, and most tissues are not strongly correlated with mortality. We developed prediction models for reproductive and single-sex tissues (Supplementary Figure S4a, S4b), but did not include them in our study of mortality due to their relatively insignificant contribution to mortality [4] [27]. During the feature selection phase of our modeling, simply using the genes with high absolute Pearson correlation of expression levels to phenotypic age yielded better model performance than using the genes identified as differentially-expressed with age by DESeq2 [25] which is a much more commonly used pipeline in bioinformatics for differential gene expression analysis [36].

Our study uncovers numerous transcriptomic biomarkers of aging in the form of model weights assigned to features, i.e., coefficients assigned to the expression values of genes. Although the reliability of many of these identified biomarkers is supported by the body of existing literature from experimental studies, a large number of them have not been explored at all. Within the list of genes corresponding to our model coefficients, there is abundant scope for future research in the field of experimental molecular biology. Such research should aim to empirically determine the correlation of the expression of specific genes with diseases of the associated tissue and ultimately mortality. As the body of human transcriptomic data grows larger, deep learning methods may also be explored in the area of transcriptomic age modeling.

## References

- [1] Maria Chiara Barbera, Laura Di Rito, Luca Guarnera, Ilaria Craparotta, Arianna Vallerga, Margherita Romeo,

- Sarah Mapelli, and Marco Bolis. Increased expression of ectodysplasin a2 receptor eda2r is the most remarkable and ubiquitous aging-related transcriptional hallmark. *Research Square*, 2021.
- [2] Reinier A Boon, Kazuma Iekushi, Stefanie Lechner, Timon Seeger, Ariane Fischer, Susanne Heydt, David Kaluza, Karine Tréguer, Guillaume Carmona, Angelika Bonauer, et al. MicroRNA-34a regulates cardiac aging and function. *Nature*, 495(7439):107–110, 2013.
  - [3] Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
  - [4] Centers for Disease Control and Prevention, National Center for Health Statistics. Leading causes of death: Data brief. Technical report, Centers for Disease Control and Prevention, 2023. Accessed: 2024-11-05. Available at: <https://www.cdc.gov/nchs/data/databriefs/db492-tables.pdf>.
  - [5] Pratishta Chatterjee, Steve Pedrini, Erik Stoops, Kathryn Goozee, Victor L Villemagne, Prita R Asih, Inge MW Verberk, Preeti Dave, Kevin Taddei, Hamid R Sohrabi, et al. Plasma glial fibrillary acidic protein is elevated in cognitively normal older adults at risk of alzheimer’s disease. *Translational psychiatry*, 11(1):27, 2021.
  - [6] Jon H Chung, Andrew R Larsen, Evan Chen, and Fred Bunz. A ptch1 homolog transcriptionally activated by p53 suppresses hedgehog signaling. *Journal of Biological Chemistry*, 289(47):33020–33031, 2014.
  - [7] GTEx Consortium. Gtex bulk tissue expression data, 2023. Accessed: 2024-11-05. Available at: [https://www.gtexportal.org/home/downloads/adult-gtex/bulk\\_tissue\\_expression](https://www.gtexportal.org/home/downloads/adult-gtex/bulk_tissue_expression).
  - [8] GTEx Consortium. Gtex metadata, 2023. Accessed: 2024-11-05. Available at: <https://www.gtexportal.org/home/downloads/adult-gtex/metadata>.
  - [9] Harris Drucker, Christopher J Burges, Linda Kaufman, Alex Smola, and Vladimir Vapnik. Support vector regression machines. *Advances in neural information processing systems*, 9, 1996.
  - [10] Sadanand Fulzele, Bharati Mendhe, Andrew Khayrullin, Maribeth Johnson, Helen Kaiser, Yutao Liu, Carlos M Isales, and Mark W Hamrick. Muscle-derived mir-34a increases with age in circulating extracellular vesicles and induces senescence of bone marrow stem cells. *Aging (Albany NY)*, 11(6):1791, 2019.
  - [11] Larry Gold, Deborah Ayers, Jennifer Bertino, Christopher Bock, Ashley Bock, Edward Brody, Jeff Carter, Virginia Cunningham, Andrew Dalby, Bruce Eaton, et al. Aptamer-based multiplexed proteomic technology for biomarker discovery. *Nature Precedings*, pages 1–1, 2010.
  - [12] Jason E Hawkes, Pamela B Cassidy, Prashiela Manga, Raymond E Boissy, David Goldgar, Lisa Cannon-Albright, Scott R Florell, and Sancy A Leachman. Report of a novel oca2 gene mutation and an investigation of oca2 variants on melanoma risk in a familial melanoma pedigree. *Journal of dermatological science*, 69(1):30–37, 2013.
  - [13] Wenjing He, Yi Yu, Wei Huang, Guoliang Feng, and Junhe Li. The pseudogene duxap8 promotes colorectal cancer cell proliferation, invasion, and migration by inducing epithelial-mesenchymal transition through interacting with ezh2 and h3k27me3. *Oncotargets and therapy*, pages 11059–11070, 2020.
  - [14] Mitchell Ho, Tapan K Bera, Mark C Willingham, Masanori Onda, Raffit Hassan, David FitzGerald, and Ira Pastan. Mesothelin expression in human lung cancer. *Clinical cancer research*, 13(5):1571–1575, 2007.
  - [15] Ina Jeong, Jae-Hyun Lim, Jin-Soo Park, and Yeon-Mok Oh. Aging-related changes in the gene expression profile of human lungs. *Aging (Albany NY)*, 12(21):21391, 2020.
  - [16] Kayla A Johnson and Arjun Krishnan. Human pan-body age-and sex-specific molecular phenomena inferred from public transcriptome data using machine learning. *bioRxiv*, pages 2023–01, 2023.
  - [17] Stefan S Kachala, Adam J Bograd, Jonathan Villena-Vargas, Kei Suzuki, Elliot L Servais, Kyuichi Kadota, Joanne Chou, Camelia S Sima, Eva Vertes, Valerie W Rusch, et al. Mesothelin overexpression is a marker of tumor aggressiveness and is associated with reduced recurrence-free and overall survival in early-stage lung adenocarcinoma. *Clinical cancer research*, 20(4):1020–1028, 2014.
  - [18] Kenneth K Kidd, Andrew J Pakstis, Michael P Donnelly, Ozlem Bulbul, Lotfi Cherni, Cemal Gurkan, Longli Kang, Hui Li, Libing Yun, Peristera Paschou, et al. The distinctive geographic patterns of common pigmentation variants at the oca2 gene. *Scientific reports*, 10(1):15433, 2020.
  - [19] Hirokazu Kimura, Alison P Klein, Ralph H Hruban, and Nicholas J Roberts. The role of inherited pathogenic cdkn2a variants in susceptibility to pancreatic cancer. *Pancreas*, 50(8):1123–1130, 2021.
  - [20] Xiqian Lan, Vinod Kumar, Alok Jha, Rukhsana Aslam, Haichao Wang, Kehong Chen, Yueming Yu, Weimei He, Feilan Chen, Huairong Luo, et al. Eda2r mediates podocyte injury in high glucose milieu. *Biochimie*, 174:74–83, 2020.
  - [21] JA Laurence and SH Fatemi. Glial fibrillary acidic protein is elevated in superior frontal, parietal and cerebellar cortices of autistic subjects. *The Cerebellum*, 4:206–210, 2005.
  - [22] Tiantian Li, Xu Li, Rongchen Mao, Lihua Pan, Yuhui Que, Chao Zhu, Lai Jin, and Shengnan Li. Nlrp2 inhibits cell proliferation and migration by regulating emt in lung adenocarcinoma cells. *Cell biology international*, 46(4):588–598, 2022.
  - [23] Xiao Ling Li, Murugan Subramanian, Matthew F Jones, Ritu Chaudhary, Deepak K Singh, Xinying Zong, Berkley Gryder, Sivasish Sindri, Min Mo, Aaron Schetter, et al. Long noncoding rna purpl suppresses basal p53 levels and promotes tumorigenicity in colorectal cancer. *Cell reports*, 20(10):2408–2423, 2017.
  - [24] John Lonsdale, Jeffrey Thomas, Mike Salvatore, Rebecca Phillips, Edmund Lo, Saboor Shad, Richard Hasz, Gary Walters, Fernando Garcia, Nancy Young, et al. The genotype-tissue expression (gtex) project. *Nature genetics*, 45(6):580–585, 2013.
  - [25] Michael I Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome biology*, 15:1–21, 2014.

- [26] Boris Muzellec, Maria Telenczuk, Vincent Cabeli, and Mathieu Andreux. Pydeseq2: a python package for bulk rna-seq differential expression analysis. *Bioinformatics*, 2023.
- [27] National Cancer Institute, Surveillance, Epidemiology, and End Results Program. Cancer stat facts: Common cancer sites, 2024. Accessed: 2024-11-05. Available at: <https://seer.cancer.gov/statfacts/html/common.html>.
- [28] National Center for Biotechnology Information. Hardy scale variable - phs000424.v4.p1, 2025. Accessed: 2025-03-01.
- [29] Hamilton Se-Hwee Oh, Jarod Rutledge, Daniel Nachun, Róbert Pálavics, Olamide Abiose, Patricia Moran-Losada, Divya Channappa, Deniz Yagmur Urey, Kate Kim, Yun Ju Sung, et al. Organ aging signatures in the plasma proteome track health and disease. *Nature*, 624(7990):164–172, 2023.
- [30] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [31] Cole T Payne, Sidra Tabassum, Silin Wu, Heng Hu, Aaron M Gusdon, Huimahn A Choi, and Xuefang S Ren. Role of microRNA-34a in blood–brain barrier permeability and mitochondrial function in ischemic stroke. *Frontiers in Cellular Neuroscience*, 17:1278334, 2023.
- [32] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [33] Angela Raucci, Federica Macrì, Stefania Castiglione, Ileana Badi, Maria Cristina Vinci, and Estella Zuccolo. MicroRNA-34a: the bad guy in age-related vascular diseases. *Cellular and Molecular Life Sciences*, pages 1–24, 2021.
- [34] Xu Ren and Pei Fen Kuan. Rnaagecalc: A multi-tissue transcriptional age calculator. *PLoS One*, 15(8):e0237006, 2020.
- [35] Rogério Ribeiro, Athos Moraes, Marta Moreno, and Pedro G Ferreira. Integration of multi-modal datasets to estimate human aging. *Machine Learning*, 113(10):7293–7317, 2024.
- [36] Diletta Rosati, Maria Palmieri, Giulia Brunelli, Andrea Morrione, Francesco Iannelli, Elisa Frullanti, and Antonio Giordano. Differential gene expression analysis pipelines and bioinformatic tools for the identification of specific biomarkers: a review. *Computational and structural biotechnology journal*, 2024.
- [37] Martina Rossi, Nirad Banksota, Chang Hoon Shin, Carlos Anerillas, Dimitrios Tsitsipatis, Jen-Hao Yang, Rachel Munk, Jennifer L Martindale, Xiaoling Yang, Yulan Piao, et al. Increased ptchd4 expression via m6a modification of ptchd4 mrna promotes senescent cell survival. *Nucleic Acids Research*, page gkae322, 2024.
- [38] S Sarkar, S Jun, S Rellick, DD Quintana, JZ Cavendish, and JW Simpkins. Expression of microRNA-34a in alzheimer’s disease brain targets genes linked to synaptic plasticity, energy metabolism, and resting state network activity. *Brain research*, 1646:139–151, 2016.
- [39] Maroun Bou Sleiman, Pooja Jha, Riekelt Houtkooper, Robert W Williams, Xu Wang, and Johan Auwerx. The gene-regulatory footprint of aging highlights conserved central regulators. *Cell reports*, 32(13), 2020.
- [40] I SomaLogic. Somascan proteomic assay technical white paper, 2015.
- [41] Rosa Spinelli, Pasqualina Florese, Luca Parrillo, Federica Zatterale, Michele Longo, Vittoria D’Esposito, Antonella Desiderio, Annika Nerstedt, Birgit Gustafson, Pietro Formisano, et al. Zmat3 hypomethylation contributes to early senescence of preadipocytes from healthy first-degree relatives of type 2 diabetics. *Aging Cell*, 21(3):e13557, 2022.
- [42] Manuj Tandon, Gina M Coudriet, Angela Criscimanna, Mairobys Socorro, Mouhanned Eliliwi, Aatur D Singhi, Zobeida Cruz-Monserrate, Peter Bailey, Michael T Lotze, Herbert Zeh, et al. Prolactin promotes fibrosis and pancreatic cancer progression. *Cancer research*, 79(20):5316–5327, 2019.
- [43] Mathias Uhlen, Linn Fagerberg, Björn M Hallström, Cecilia Lindskog, Per Oksvold, Adil Mardinoglu, Åsa Sivertsson, Caroline Kampf, Evelina Sjöstedt, Anna Asplund, et al. Tissue-based map of the human proteome. *Science*, 347(6220):1260419, 2015.
- [44] Jingyun Wang, Yang Mi, Xiangdong Sun, Xia Xue, Huanjie Zhao, Mengfei Zhang, Baitong Hu, Ihtisham Bukhari, and Pengyuan Zheng. Lnc-ptchd4-as inhibits gastric cancer through msh2-msh6 dimerization and atm-p53-p21 activation. *Aging (Albany NY)*, 15(22):13558, 2023.
- [45] Ping Wang, Zhou Fang, Wei Pei, Qi Wu, Tingting Niu, Chengyuan Dong, Mingkang Wu, Bei Li, and Zhijie Gao. Senescence reprogramming by mthfd2 deficiency facilitates tumor progression. *J Cancer*, 2024.
- [46] Svante Wold, Michael Sjöström, and Lennart Eriksson. Pls-regression: a basic tool of chemometrics. *Chemometrics and intelligent laboratory systems*, 58(2):109–130, 2001.
- [47] Shengkai Xia, Wenzhe Duan, Mingxin Xu, Mengqi Li, Mengyi Tang, Song Wei, Manqing Lin, Encheng Li, Wenwen Liu, and Qi Wang. Mesothelin promotes brain metastasis of non-small cell lung cancer by activating met. *Journal of Experimental & Clinical Cancer Research*, 43(1):103, 2024.
- [48] Lin Zhu, Xianhui Liu, Weiyu Zhang, Hao Hu, Qi Wang, and Kexin Xu. Mthfd2 is a potential oncogene for its strong association with poor prognosis and high level of immune infiltrates in urothelial carcinomas of bladder. *BMC cancer*, 22(1):556, 2022.

# 1 Supplementary Materials

## 1.1 The Hardy Scale

The Hardy Scale is a categorical measure used to assess the overall health status of individuals at the time of death, in the Genotype-Tissue Expression Project (GTEx). It provides a standardized framework for classifying subjects based on the presence and severity of disease conditions. The scale is typically assigned based on clinical and pathological evaluations and is often used in research involving aging, mortality risk, and disease progression. The scale consists of the following categories:

- **Type 1 (Violent and fast death of unnatural causes):** Deaths due to accident, blunt force trauma, or suicide, terminal phase estimated at less than 10 minutes.
- **Type 2 (Fast death of natural causes):** Sudden unexpected deaths of people who had been reasonably healthy, after a terminal phase estimated at less than 1 hour (with sudden death from a myocardial infarction as a model cause of death for this category).
- **Type 3 (Intermediate death):** Death after a terminal phase of 1 to 24 hours (not classifiable as 2 or 4); patients who were ill but death was unexpected.
- **Type 4 (Slow death):** Death after a long illness, with a terminal phase longer than 1 day (commonly cancer or chronic pulmonary disease); deaths that are not unexpected.
- **Type 0 (Ventilator case):** All cases on a ventilator immediately before death.

In studies on the GTEx dataset, including this one, the Hardy Scale ratings of subjects are typically defined by the variable named *dthrdy*, as done in the GTEx study.

## 1.2 Age range distribution in GTEx phenotypes

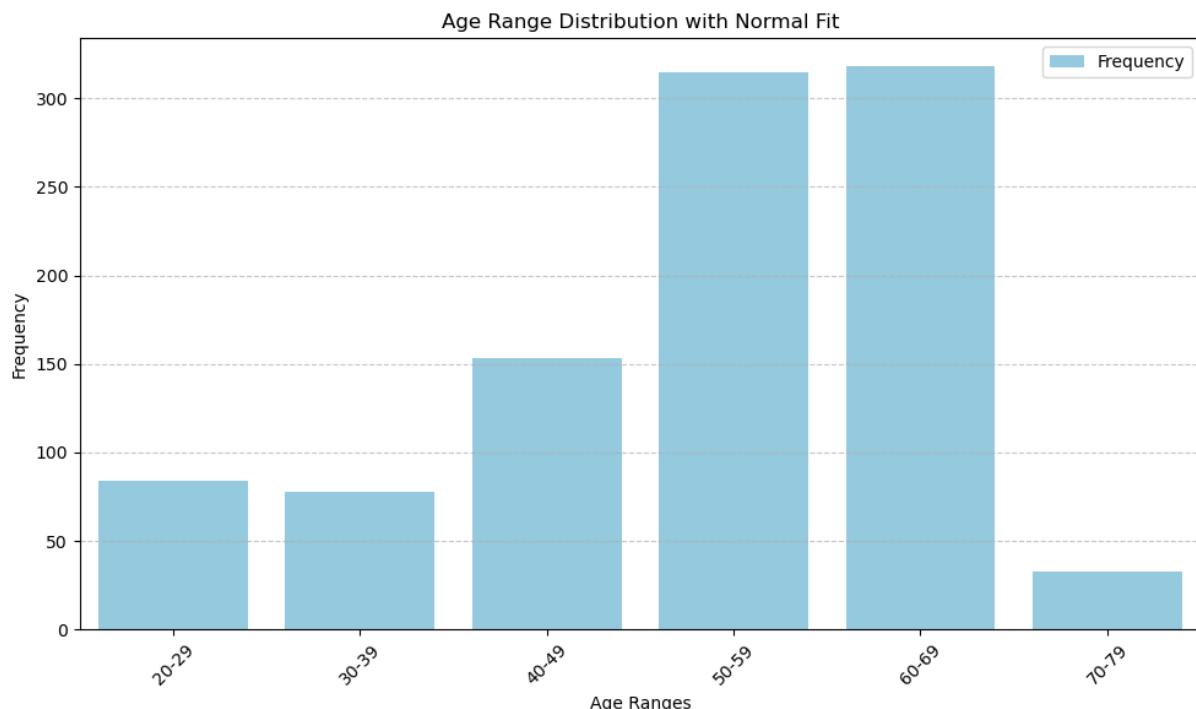
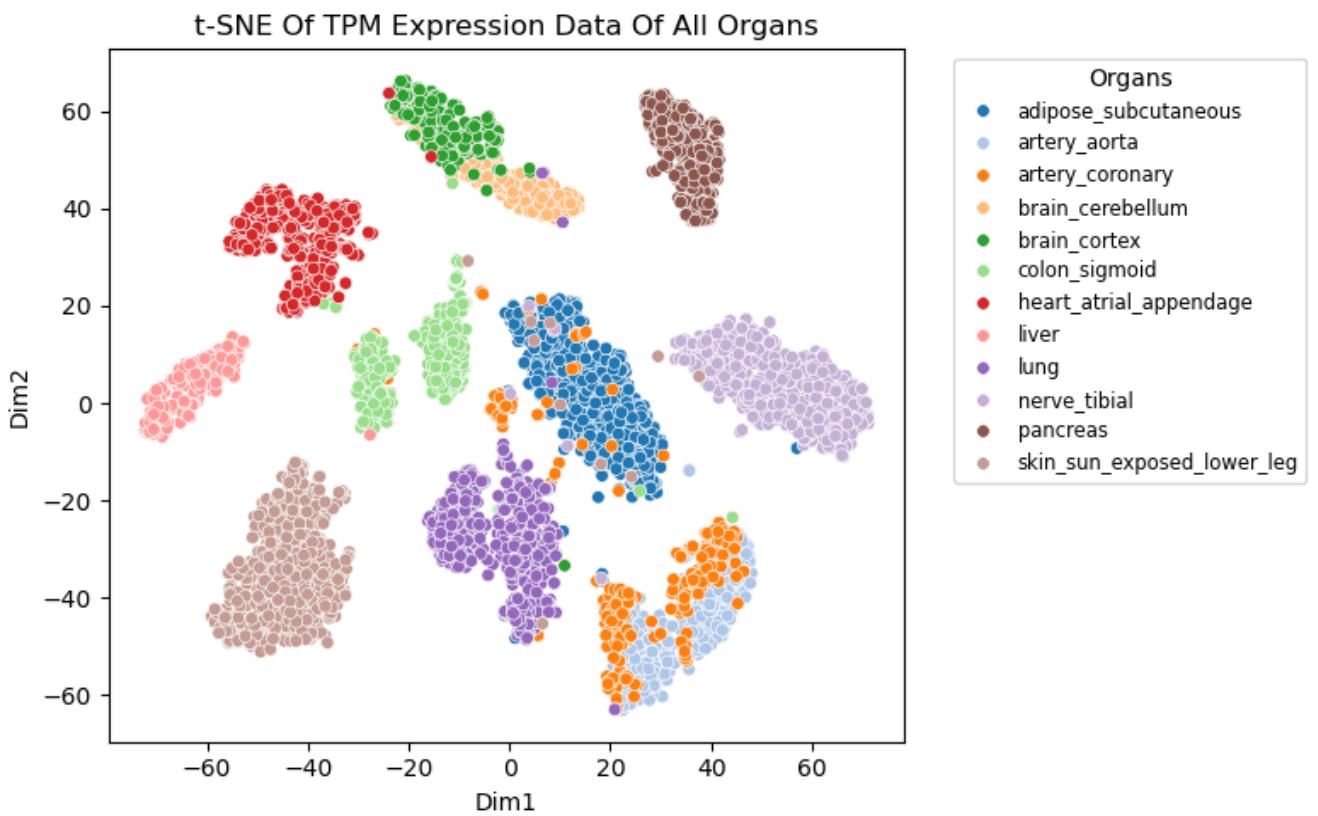


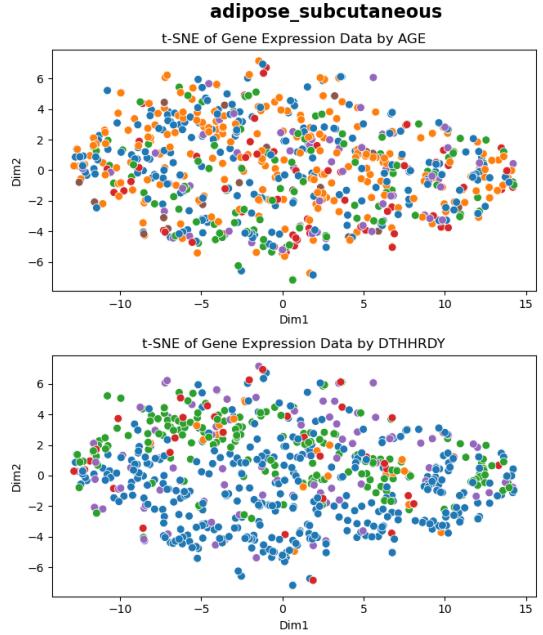
Figure S1: Age distribution of GTEx subjects.

### 1.3 Clustering of gene expression (TPM) data of the samples from 12 types of tissues

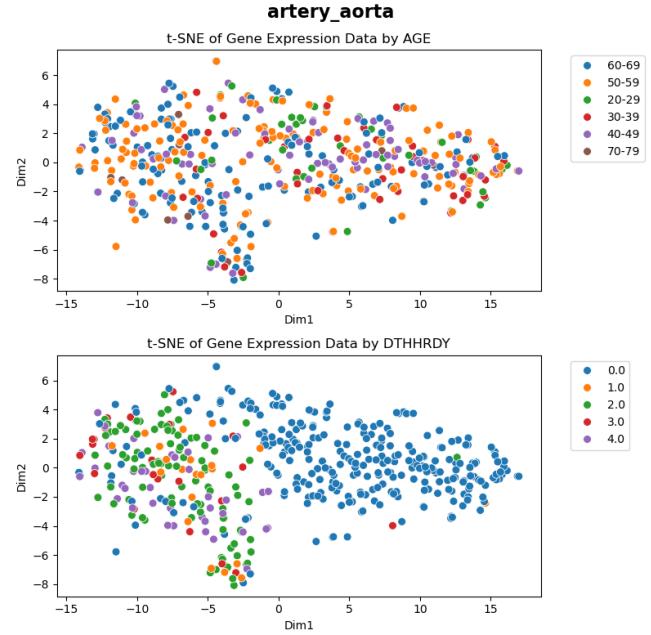


**Figure S2: T-SNE plot of the tissue-specific gene expression (TPM) values.** The two axes, Dim1 and Dim2, are low-dimensional embeddings learned by the t-SNE algorithm to best preserve local structure and relationships between high-dimensional gene expression profiles. Dots represent the RNA-Seq samples of the 12 tissue types from the Adult GTEx that we have used in our study. It is observed that the samples of the same type of tissue strongly tend to cluster together with themselves while remaining noticeably separate from the samples of other tissue types. It is also evident that tissues from the same or related organs, such as the coronary artery and aorta, or, the brain's cerebellum and cortex, cluster adjacent to each other.

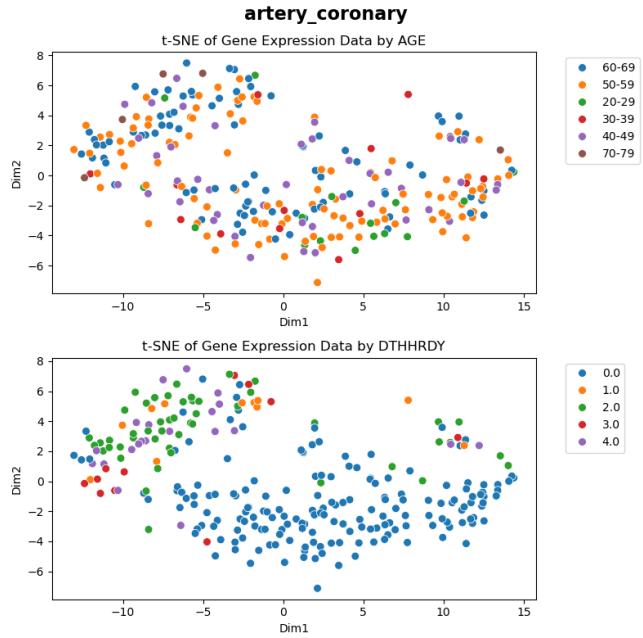
## 1.4 Clustering of tissue-specific gene expression (TPM) data



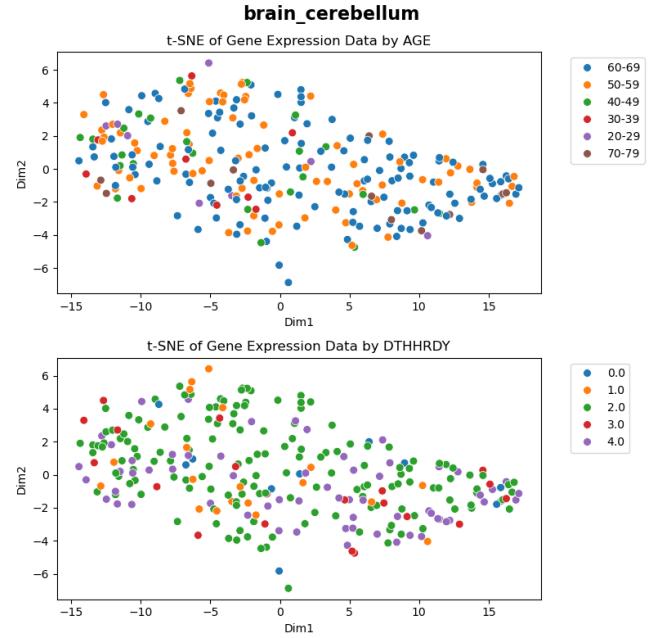
(a) Clustering of subcutaneous adipose tissue TPM expression data



(b) Clustering of aorta tissue TPM expression data

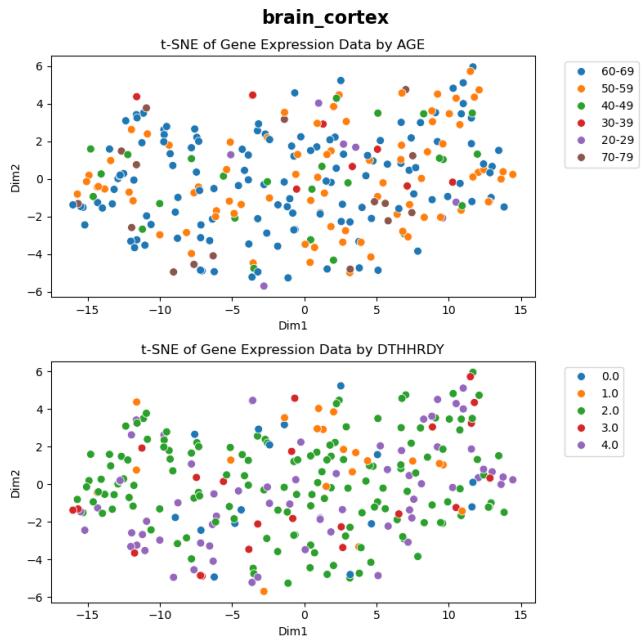


(c) Clustering of coronary artery tissue TPM expression data

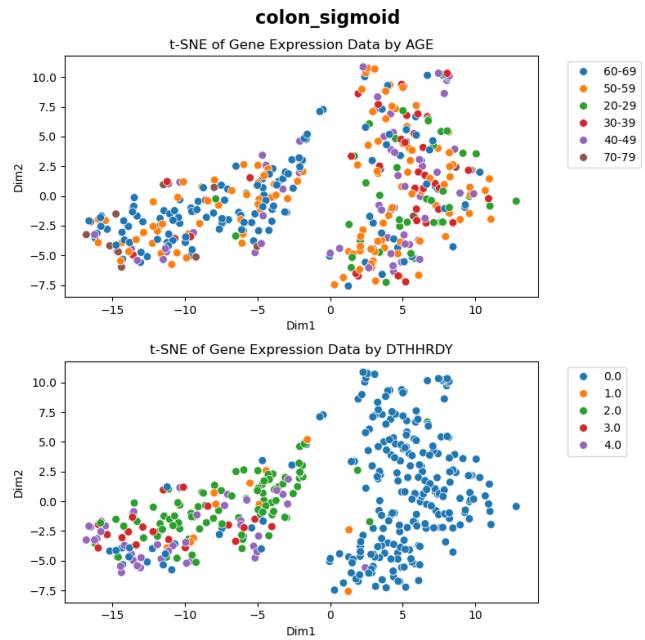


(d) Clustering of brain cerebellum tissue TPM expression data

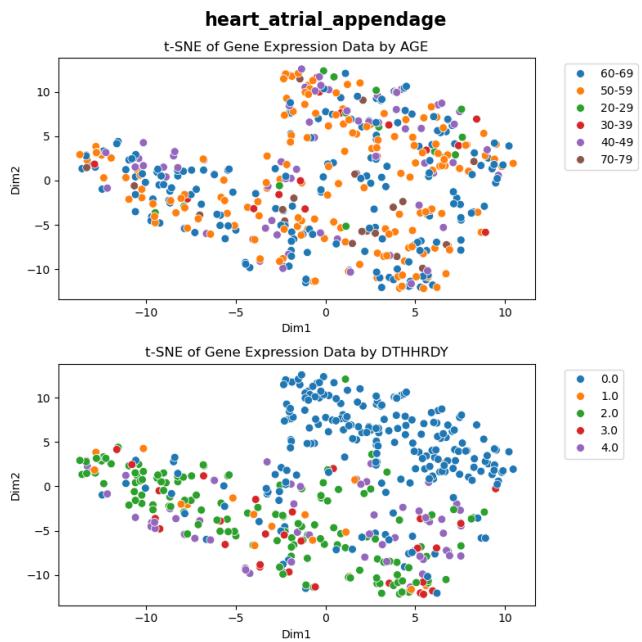
**Figure S3: Part 1: t-SNE plots of gene expression (TPM) for individual tissues (1–4).** In most tissues, it is observed that the samples from subjects with a Hardy Scale rating of 0 (death on a ventilator) tend to cluster together with themselves and separate from the samples from subjects with other death types. No significant grouped clustering by age range can be observed.



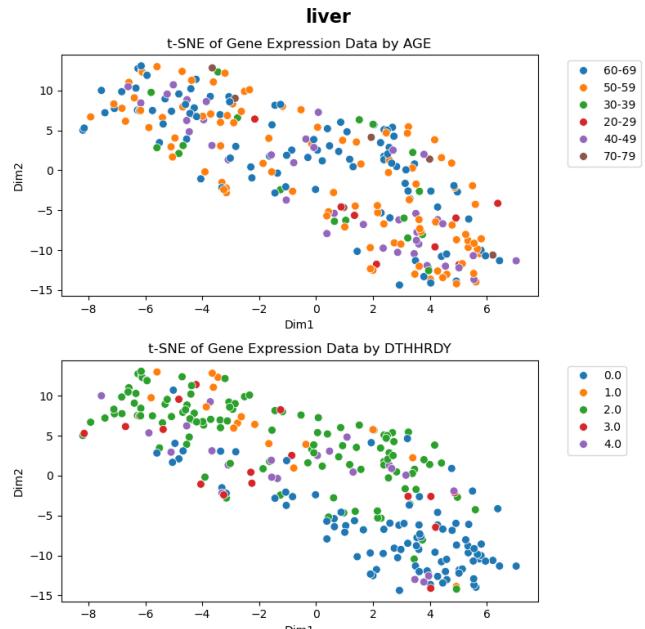
(e) Clustering of brain cortex tissue TPM expression data



(f) Clustering of sigmoid colon tissue TPM expression data

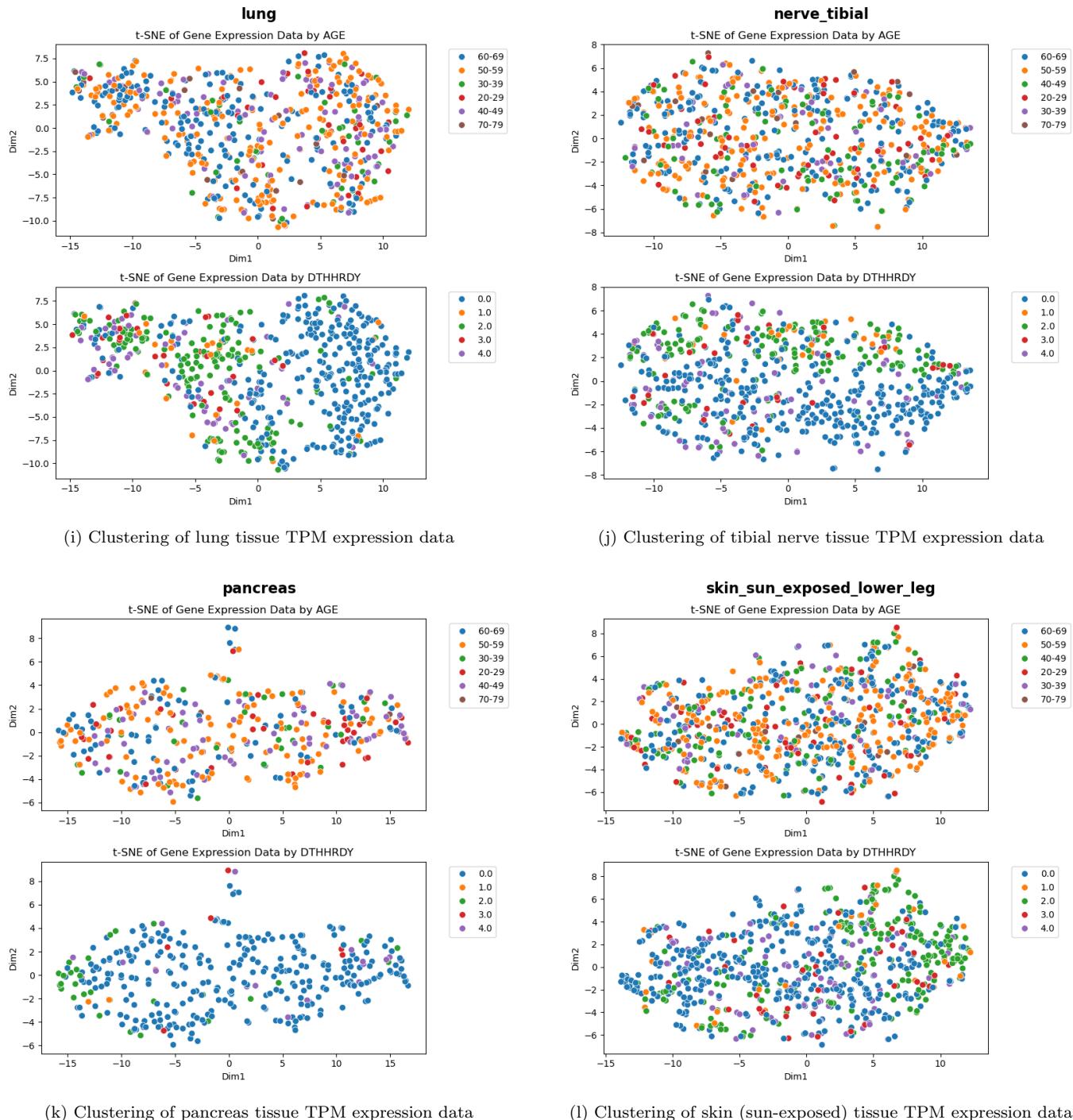


(g) Clustering of heart atrial appendage tissue TPM expression data



(h) Clustering of liver tissue TPM expression data

Figure S3: Part 2: t-SNE plots of gene expression (TPM) for individual tissues (1–4). In most tissues, it is observed that the samples from subjects with a Hardy Scale rating of 0 (death on a ventilator) tend to cluster together with themselves and separate from the samples from subjects with other death types. No significant grouped clustering by age range can be observed.



**Figure S3: Part 3: t-SNE plots of gene expression (TPM) for individual tissues (1–4).** In most tissues, it is observed that the samples from subjects with a Hardy Scale rating of 0 (death on a ventilator) tend to cluster together with themselves and separate from the samples from subjects with other death types. No significant grouped clustering by age range can be observed.

## 1.5 Conditional probability analysis of maximum age-gap of all subjects

DTHHRDY	$\text{Agegap} \geq 1.5353$	$\text{Agegap} \leq 0.8033$	$1.0208 < \text{Agegap} < 1.3148$	Overall ( $p$ )
0	0.3377	0.6980	0.3864	0.5066
1	0.0265	0.0336	0.0568	0.0377
2	0.3974	0.1879	0.3636	0.2994
3	0.0662	0.0201	0.0682	0.0508
4	0.1722	0.0604	0.1250	0.1055

Table 1: Conditional probabilities of DTHHRDY given different age-gap ranges

The following observations were made on the basis of the calculated conditional probabilities:

- **$dthhrdy=0$ :**
  - Extreme negative agers are 2.067 times as likely to have died with  $dthhrdy=0$  compared to extreme positive agers.
  - Average agers are 1.144 times as likely to have died with  $dthhrdy=0$  compared to extreme positive agers.
  - Extreme negative agers are 1.807 times as likely to have died with  $dthhrdy=0$  compared to average agers.
- **$dthhrdy=1$ :**
  - Extreme negative agers are 1.267 times as likely to have died with  $dthhrdy=1$  compared to extreme positive agers.
  - Average agers are 2.145 times as likely to have died with  $dthhrdy=1$  compared to extreme positive agers.
  - Average agers are 1.693 times as likely to have died with  $dthhrdy=1$  compared to extreme negative agers.
- **$dthhrdy=2$ :**
  - Extreme positive agers are 2.114 times as likely to have died with  $dthhrdy=2$  compared to extreme negative agers.
  - Extreme positive agers are 1.093 times as likely to have died with  $dthhrdy=2$  compared to average agers.
  - Average agers are 1.935 times as likely to have died with  $dthhrdy=2$  compared to extreme negative agers.
- **$dthhrdy=3$  or  $4$ :**
  - Extreme positive agers are 2.96 times as likely to have died with  $dthhrdy=3$  or 4 compared to extreme negative agers.
  - Extreme positive agers are 1.234 times as likely to have died with  $dthhrdy=3$  or 4 compared to average agers.
  - Average agers are 2.399 times as likely to have died with  $dthhrdy=3$  or 4 compared to extreme negative agers.

## 1.6 Performance comparison of pipeline variants

Tissue Type	Metric	Elastic Net		PLS		SVR		Rand. Forest		Lasso
		Corr.	DEG	Corr.	DEG	Corr.	DEG	Corr.	DEG	Enriched
Liver	RMSE	6.02	6.78	6.41	6.81	6.41	6.86	8.14	8.32	8.75
	$R^2$	0.54	0.42	0.48	0.42	0.48	0.41	0.17	0.13	0.04
Aorta	RMSE	5.32	5.62	5.48	6.01	5.51	6.02	6.43	6.00	8.91
	$R^2$	0.74	0.71	0.73	0.67	0.72	0.67	0.62	0.67	0.28
Coronary Artery	RMSE	5.39	6.81	5.40	7.06	5.47	7.19	8.23	8.65	10.27
	$R^2$	0.73	0.57	0.73	0.54	0.73	0.53	0.38	0.31	0.03
Brain Cortex	RMSE	6.46	6.23	6.43	6.06	6.47	6.11	8.42	8.10	7.93
	$R^2$	0.59	0.62	0.60	0.64	0.59	0.64	0.31	0.36	0.39
Brain Cerebellum	RMSE	7.88	7.08	8.36	7.21	8.39	6.94	8.44	8.02	9.67
	$R^2$	0.43	0.54	0.36	0.52	0.35	0.56	0.34	0.41	0.14
Heart Atr. App.	RMSE	7.23	6.44	6.99	6.48	6.63	6.49	8.33	8.43	8.63
	$R^2$	0.36	0.49	0.40	0.48	0.46	0.48	0.15	0.13	0.09
Subcut. Adipose	RMSE	7.01	7.57	7.38	7.89	7.42	7.91	8.66	9.25	10.31
	$R^2$	0.62	0.56	0.58	0.52	0.58	0.52	0.42	0.34	0.18
Lung	RMSE	6.27	7.46	6.41	7.49	6.41	7.80	9.03	8.99	9.55
	$R^2$	0.66	0.52	0.65	0.52	0.65	0.48	0.30	0.31	0.22
Sun-Exposed Skin	RMSE	6.81	6.12	6.76	6.69	7.45	6.70	8.73	9.02	9.57
	$R^2$	0.68	0.74	0.68	0.69	0.61	0.69	0.47	0.43	0.36
Tibial Nerve	RMSE	5.44	6.51	5.69	6.50	5.69	6.51	7.51	7.95	9.26
	$R^2$	0.79	0.70	0.77	0.70	0.77	0.70	0.60	0.55	0.39
Sigmoid Colon	RMSE	6.52	7.45	7.68	8.00	7.70	8.00	8.44	8.48	10.90
	$R^2$	0.67	0.57	0.54	0.50	0.54	0.50	0.44	0.44	0.07
Pancreas	RMSE	6.45	5.88	7.20	6.39	7.15	6.38	6.59	8.35	10.84
	$R^2$	0.61	0.67	0.51	0.61	0.52	0.62	0.59	0.34	-0.11

Table 2: Performance comparison of pipeline variants

## 1.7 Training time comparison of learning algorithms

Tissue Type	Lasso	Elastic Net	SVR	PLS	Random Forest
Liver	30	48	41	4	2
Sigmoid Colon	166	637	153	12	7
Sun-Exposed Skin	107	56	127	0	3
Subcutaneous Adipose	34	1526	237	27	10
Coronary Artery	8	11	16	0	1
Aorta	9	199	250	41	9
Brain Cortex	30	50	47	1	3
Brain Cerebellum	24	300	122	3	8
Heart Atrial Appendage	70	57	26	0	0
Lung	126	3151	145	2	6
Tibial Nerve	137	91	415	5	12
Pancreas	28	2	16	0	0
<b>Total</b>	<b>769</b>	<b>6128</b>	<b>1595</b>	<b>95</b>	<b>61</b>

\* All predictors except those using random forest regression were 20× bootstrapped.

Table 3: Training time (in seconds) comparison of learning algorithms.

## 1.8 Performance comparison with other studies

Study	Tissue/Organ	Model	MAE	RMSE	$R^2$
Ribeiro et al. [35]	Lung	EN*	6.29	8.33	0.59
		GBT*	7.07	8.64	0.56
	Ours*	EN*	<b>3.71</b>	<b>5.83</b>	<b>0.63</b>
		PLS	4.85	6.11	0.59
Ren et al. [34]	Liver	RAC*	9.49	11.71	-0.2
		PLS	<b>5.39</b>	<b>6.94</b>	<b>0.31</b>
Ren et al. [34]	Tibial Nerve	RAC*	10.61	12.59	0.08
		PLS	<b>4.71</b>	<b>5.94</b>	<b>0.70</b>

\* “Ours” refers to methods proposed in this study.

\* EN, GBT, and RAC refer to the Elastic net, Gradient Boosted Trees, and RNAAgeCalc [34] models respectively.

Table 4: Performance comparison with other studies.

## 1.9 Prediction models for reproductive and single-sex tissues

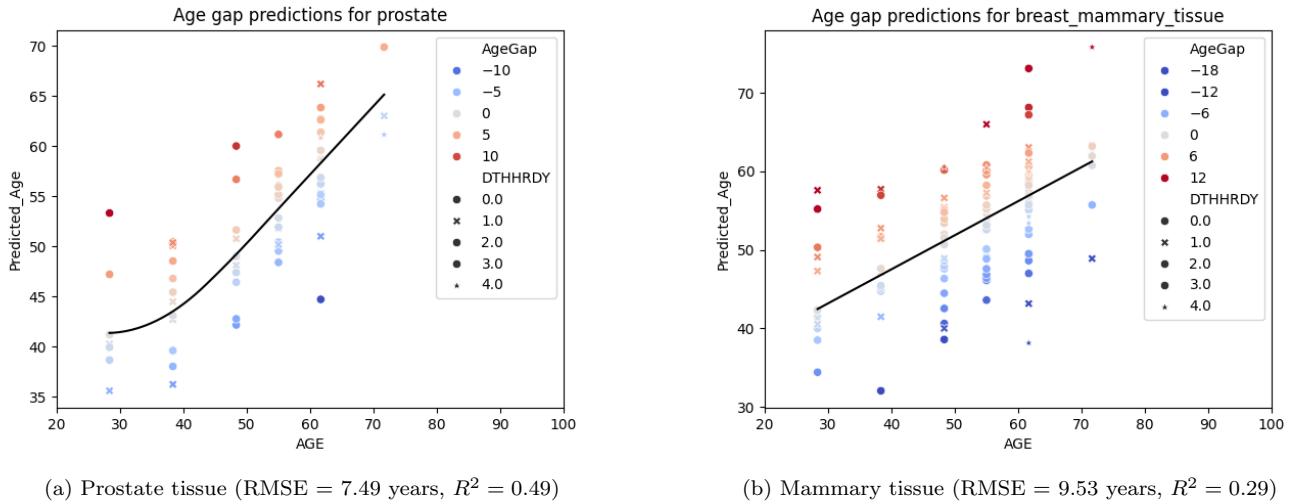


Figure S4: Age prediction performance in single-sex tissues.