

Data Analysis Project

Cross Species Queries of Large Gene Expression Databases *

Hai-Son Le

Joint work with: Ziv Bar-Joseph, and Zoltán N. Oltvai

DAP Committee: William Cohen, Ziv Bar-Joseph, and Zoltán N. Oltvai

Abstract

Motivation: Expression databases, including the Gene Expression Omnibus (GEO) and ArrayExpress, have experienced significant growth over the last decade and now hold hundreds of thousands of arrays from multiple species. Since most drugs are initially tested on model organisms, the ability to compare expression experiments across species may help identify pathways that are activated in a similar way in humans and other organisms. However, while several methods exist for finding co-expressed genes in the same species as a query gene, looking at co-expression of homologs or arbitrary genes in other species is challenging. Unlike sequence, which is static, expression is dynamic and changes between tissues, conditions and time. Thus, to carry out cross species analysis using these databases we need methods that can match experiments in one species with experiments in another species.

Results: To facilitate queries in large databases we developed a new method for comparing expression experiments from different species. We define a distance metric between the ranking of orthologous genes in the two species. We show how to solve an optimization problem for learning the parameters of this function using a training dataset of known similar expression experiments pairs. The function we learn outperforms previous methods and simpler rank comparison methods that have been used in the past for single species analysis. We used our method to compare millions of array pairs from mouse and human expression experiments. The resulting matches can be used to find functionally related genes, to hypothesize about biological response mechanisms and to highlight conditions and diseases that are activating similar pathways in both species.

Availability: Supporting methods, results are available from <http://sb.cs.cmu.edu/ExpQ/>.

1 Introduction

Advances in sequencing technology have led to a remarkable growth in the size of sequence databases over the last two decades. This has allowed researchers to study newly sequenced genes by utilizing knowledge about their homologs in other species [16]. Alignment and search methods, most notably BLAST [1], have become standard tools and are extensively used by molecular biologists. Cross species analysis of sequence data is now a standard practice. However, similar usage of expression databases has not materialized. Expression databases, including Gene Expression Omnibus (GEO; www.ncbi.nih.gov/geo/) and ArrayExpress (www.ebi.ac.uk/Databases/microarray.html) hold hundreds of thousands of arrays from multiple species (see Fig. 1). Co-expression is a powerful method for assigning new function to genes within a single species [21]. If we are able to identify a large set of matched expression experiments across species this method can be extended and used in a cross species analysis setting as well. Consider a human gene with unknown function that is co-expressed (across many different conditions) with a mouse gene with known function. This information can provide useful clues about the function of the human gene. This information is also

*Earlier version appears in [15].

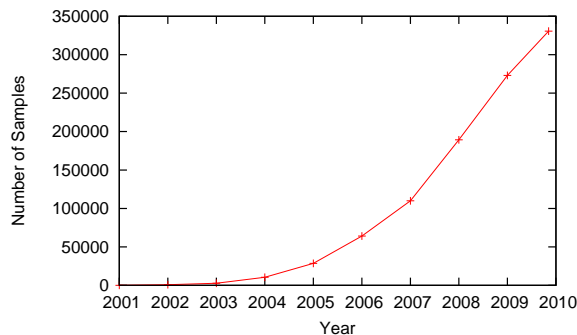


Figure 1: Growth of microarray databases. Growth in microarray datasets deposited in GEO in the last decade. The growth resembles the impressive growth of sequence databases in the 90’s.

useful for identifying orthologs. If a gene has multiple homologs in another species then the homolog with the highest co-expression similarity in several conditions is likely its orthologs since they are involved in the same processes in both species.

While promising, querying expression datasets to identify co-expressed genes in other species is challenging. Unlike sequence, which is static, expression is dynamic and changes between tissues, conditions and time. Thus, a key challenge is to match experiments in one species with experiments in another species. Almost all studies that have analyzed expression datasets in multiple species relied on one of two methods. They have either carried out experiments under the same condition in multiple species or have looked at co-expression within a species and tested whether these relationships are retained across species. Examples of the former set of methods include comparison of cell cycle experiments across species [14], comparing response programs [17] and comparing tissue expression between human and mouse [24]. Examples of the latter strategy include the metaGene analysis [23] and cross species clustering methods [18]. See [19] for a recent review of these methods.

While successful, the approaches discussed above are not appropriate for querying large databases. In almost all cases it is impossible to find a perfect match for a specific condition in the database. Even in the rare cases when such matches occur it is not clear if the same pathways are activated in the different species. For example, many drugs that work well on animal models fail when applied to humans, at least in part because of differences in the pathways involved [5]. Looking at relationships within and between species would also not answer the questions we mentioned above since these require knowledge of orthologs assignment to begin with. These methods are also less appropriate for identifying one to one gene matchings because they are focusing on clusters instead.

The only previous attempt we are aware of to facilitate cross species queries of expression data is the non-negative matrix factorization (NMF) approach presented by Tamayo et al [25]. This unsupervised approach discovers a small number of metagenes (similar to principle components) that capture the invariant biological features of the dataset. The orthologs of the genes included in the metagenes are then combined in a similar way in the query species to identify related expression datasets. While the approach was successfully used to compare two specific experiments in humans and mouse, as we show in Results, the fact that the approach is unsupervised makes it less appropriate for large scale queries of expression databases.

In this paper, we present a new method for identifying similar experiments in different species. Instead of relying on the description of the experiments we develop a method to determine the similarity of expression profiles by introducing a new distance function and utilizing a group of known orthologs. Our method uses a training dataset of known similar pairs to learn the parameters for distance functions between pairs of experiments based on the rank of orthologous genes overcoming problems related to difference in noise and platforms between species. We show that the function we learn outperforms simpler rank comparison methods that have been used in the past [10, 13]. We next use our method to compare millions of array pairs from mouse and human experiments. The resulting matches highlight conditions and diseases that

are activating similar pathways in both species and can also hint at diseases were these pathways seem to differ. Given the large number of arrays in current databases our methods can also be used to aid manual annotations of cross species similarity by focusing on a small subset of the millions of possible matches.

We note that while the discussion below focuses on microarray data and we have only tested our methods on such data, our methods are appropriate for deep sequencing expression data as well. As long as a partial orthologs list can be obtained the methods we present below can be used to compare any expression datasets across species.

2 Methods

2.1 Comparing microarrays across species

Our goal is to obtain a distance function that given two microarray datasets outputs a small distance between experiments that are very similar and a large distance for those pairs that study different processes or in which different pathways are activated in the two species being compared. Since we are comparing experiments from different platforms and species the first decision we made was to compare the ranking of the genes in each array rather than their expression levels (previous methods for comparing experiments in the same species have relied on ranking as well [10]). There are a number of other properties that we seek for such scoring functions. First, they should of course be able to separate similar pairs from non similar pairs. In addition, it would be useful if the function is a metric or a pseudometric (a pseudometric satisfies all properties of a metric except for the identity, that is $d(x, y)$ could be 0 even if $x \neq y$). This will guarantee useful distance properties including symmetry and triangle inequality (see Supporting Methods for the complete list). Finally, we would like to be able to determine some statistical properties for these scoring methods in order to determine a p-value for the similarity / difference between the experiments being compared (Section 2.3.1).

2.1.1 Notations

We first provide notations that are used in the rest of the paper. As mentioned above our function would be constructed from metrics on permutations (ordering) of ranks. Each microarray experiment is a vector in R^n , where each dimension is the expression value for a specific gene. We consider the problem of comparing a microarray \mathcal{X} of a species A with n_A genes and a microarray \mathcal{Y} of a species B with n_B genes. There are m orthologs between the two species. In other words, there is a one-to-one mapping O from m species A genes to m species B genes. $1, \dots, m$ are the orthologs, $X = \{X_i : 1 \leq i \leq m\}$ and $Y = \{Y_i : 1 \leq i \leq m\}$ are the expression values of the orthologs in \mathcal{X} and \mathcal{Y} , respectively. Let π, σ be the rank orderings of the expression values of the orthologs in X and Y . For simplicity, we assume that there are no ties in rankings. Therefore, π, σ are two elements of the permutation group G_m . Recall that $\pi, \sigma : \{1, \dots, m\} \rightarrow \{1, \dots, m\}$ are bijections: $\pi(i), \sigma(i)$ are the ranks given to the ortholog i , with lowered numbered ranks given to higher expression values. Also let I_m be the identity permutation in G_m . Finally, $\text{tr}(M)$ is the trace of a matrix M .

Assume we have a metric d on G_m . For our significance analysis we test the null hypothesis H_0 that π and σ are not associated versus the alternate hypothesis that they are. One way is to ask how large $d(\pi, \sigma)$ would be if σ were chosen uniformly at random. More formally, let D_d be the distribution of $d(\pi, \sigma)$ when σ is drawn uniformly from G_m . We reject the null hypothesis H_0 if $d(\pi, \sigma)$ is significantly smaller than $E(D_d)$. This setting is a standard approach in literature [8] (see also Supporting Fig. 1).

2.2 Fixed distance function: Spearman’s rank correlation

Below we discuss distance functions that satisfy the requirements mentioned above for cross species analysis. We first discuss a method that does not require any parameter tuning. Such methods have been extensively used for comparing permutations. However, as we show in Sect. 3.2 they are less appropriate for gene

expression data due to the unique properties of such data. In the next section we discuss modification of these methods that are more appropriate for the expression data we are working with.

The Spearman's rank correlation R metric is defined as:

$$R(\pi, \sigma) = \sqrt{\sum_{i=1}^m (\pi(i) - \sigma(i))^2} \quad (1)$$

In other words it is the L_2 distance between π and σ . Hence, it is a metric. Moreover, using Hoeffding's central limit theorem it can be proved that R^2 has a limiting normal distribution [8]. Note that frequently, R is standardized to have values in $[-1, 1]$. This yields the widely used Spearman's rank correlation ρ .

$$\rho = 1 - \frac{6R^2(\pi, \sigma)}{(m^3 - m)} \quad (2)$$

2.3 Adaptive Metrics

While fixed methods that do not require parameter tuning have proven useful for many cases they are less appropriate for expression data. In such data the importance of the ranking is not uniform. In other words genes that are expressed at very high or very low levels compared to baseline may be very informative whereas the exact ranking of genes that are expressed at baseline levels may be much less important. Thus, rank differences for genes in the middle of the rankings are more likely due to noise. An appropriate way to weight the differences between the rankings may lead to a better distance function between arrays. The key challenge is to determine what are the important ranks and how they should be weighted. Below we present a number of adaptive methods that can address this issue. The methods we present differ in the number of parameters that needs to be learned and thus each may be appropriate for different cases depending on the amount of training data that exists.

2.3.1 Weighted Rank Metric

Using a weight vector w of length m , we can modify the Spearman's rank correlation and define the following metric:

$$d(\pi, \sigma) = \sqrt{\sum_{i=1}^m (w(\pi(i)) - w(\sigma(i)))^2} \quad (3)$$

The vector w defines the weight of each rank and thus captures the significance of each rank in measuring the association of two microarrays. Consider two arrays $(1, 2, 3, 4)$ and $(1, 3, 2, 4)$. Their Spearman R distance is $\sqrt{2}$ while for a weight vector $w = (1, 0, 0, 1)$, their distance would be 0. Such a weight vector places the weight on the top and bottom matches and disregards middle orderings.

The resulting function is no longer a metric, but rather a pseudo-metric in the original π, σ space ($d(\pi, \sigma) = 0$ does not imply $\pi = \sigma$). However, it is easy to see that it is a metric in the transformed $w(\cdot)$ -space because it is a L_2 distance between the vectors $w(\pi)$ and $w(\sigma)$, where $w(\pi) = (w(\pi_1), \dots, w(\pi_m))$ and similarly for $w(\sigma)$. In other words the w -transformation makes some of the permutations indistinguishable indicating that the changes made are not significant and so the two permutations result in the same weighted vector. However, for those permutations that are still distinguishable following the w -transformation the metric properties are preserved. The distribution D_d of $d(\pi, \sigma)$ when σ is drawn uniformly from G_m is asymptotically normal. See Supporting Methods for proof. We can calculate the mean and variance of D_d through exact calculation or random sampling. P-value can then be calculated based on this normal distribution.

A specific assignment of weights which is in line with our assumptions regarding the importance of genes expression ranks is the following modified Spearman's rank correlation.

2.3.2 Top-Bottom R (TBR)

For any $0 < k < 1$ and $r > 0$ we can define w as following:

$$w(i) = \begin{cases} r(i - km) & \text{if } 1 \leq i < km, \\ r(i - (1 - k)m) & \text{if } (1 - k)m < i \leq m, \\ 0 & \text{otherwise.} \end{cases} \quad (4)$$

Note that genes expressed at a high level will have negative weights and those with low levels positive weights allowing the method to penalize experiments in which genes move from one extreme to the other. All middle ranks $[km, (1 - k)m]$ are assigned the same weight so genes that have ranks changed within this interval do not affect the distance at all. At the same time, it scales the high and low ranks r times to a wider range to increase the granularity of rank difference. Choosing the value of k and r can either be done using cross validation or it could be manually specified.

2.3.3 Learning a complete weight vector w

While the above method leads to different weights for different rankings it specifies a very strict cutoff which may not accurately represent the importance of the differences in ranking. An alternative approach is to assign weights that are continuously changing based on the ranking by learning a weight vector from training data. Here we assume that we have access to such training data which is indeed the case for a number of pairs of species (most notably tissue data for human and mouse as we use in Results). Assume we have M microarrays of species A and N microarrays of species B and for each microarray, let \mathcal{S} be the set of pairs of similar arrays and \mathcal{D} is the set of pairs of dissimilar arrays. If the dissimilar arrays are not known, we can select \mathcal{D} as the set of all pairs that are not in \mathcal{S} .

Each permutation π can be represented as a binary $m \times m$ matrix M_π .

$$M_\pi(i, j) = \begin{cases} 1 & \text{if } \pi(i) = j, \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

Using this notation we can define an L_2 metric d as:

$$d(\pi, \sigma) = \|M_\pi w - M_\sigma w\|_2 \quad (6)$$

$$= \sqrt{w^T (M_\pi - M_\sigma)^T (M_\pi - M_\sigma) w} \quad (7)$$

Our goal is to learn a vector w such that this distance be small for the positive set and large for the negative set. This leads to the following optimization problem:

$$\min \sum_{(x,y) \in \mathcal{S}} w^T (M_{\pi_x} - M_{\pi_y})^T (M_{\pi_x} - M_{\pi_y}) w \quad (8)$$

$$\text{s.t } \sum_{(x,y) \in \mathcal{D}} w^T (M_{\pi_x} - M_{\pi_y})^T (M_{\pi_x} - M_{\pi_y}) w = 1 \quad (9)$$

Note that the summation is on different groups. The optimization (top) is summed over the similar pairs whereas the constraint (bottom) is summed over the dissimilar pair. The choice of the constant 1 on the right hand side of (9) is arbitrary. However, replacing it with any constant $c > 0$ results only in w being multiplied by \sqrt{c} which leads to the same order of scores for microarray pairs and so does not change our results. We can further simplify the problem to

$$\min w^T Z_{\mathcal{S}} w \quad (10)$$

$$\text{s.t } w^T Z_{\mathcal{D}} w = 1 \quad (11)$$

with

$$Z_S = \sum_{(x,y) \in \mathcal{S}} (M_{\pi_x} - M_{\pi_y})^T (M_{\pi_x} - M_{\pi_y})$$

$$Z_D = \sum_{(x,y) \in \mathcal{D}} (M_{\pi_x} - M_{\pi_y})^T (M_{\pi_x} - M_{\pi_y})$$

The matrices Z_S and Z_D are positive semidefinite since they are sums of positive semidefinite matrices $(M_{\pi_x} - M_{\pi_y})^T (M_{\pi_x} - M_{\pi_y})$. Although this optimization is not convex, there exists global minima based on the reformulation of this problem to finding eigenvalues of the Rayleigh quotient. The derivation is similar to Fisher's Linear Discriminant Analysis [11].

2.3.4 Relational Weighted Rank Metric

A drawback of the weight vector distance metric discussed above is that it assigns weights to ranks in each microarray independent of the ranks in the other microarray. To overcome this problem we extend the vector weight w into a full matrix W to incorporate the dependence between ranks in two microarrays. For a pair of microarrays with ortholog rankings π and σ , define a symmetric $m \times m$ matrix $M_{\pi,\sigma}^F$, whose entries (i, j) are non-zeros if and only if there exists a gene g such that g is ranked i and j in the microarrays, respectively. Formally,

$$M_{\pi,\sigma}^F(i, j) = \mathbf{1} [\pi^{-1}(i) = \sigma^{-1}(j)] + \mathbf{1} [\pi^{-1}(j) = \sigma^{-1}(i)] \quad (12)$$

In other words, $M_{\pi,\sigma}^F$ is a matrix where an entry of 1 in location (i, j) indicates that the gene in location i in the first experiment is the same as the gene in location j in the second or vice versa. By definition, $M_{\pi,\sigma}^F$ is a symmetric matrix. Note that this definition implies that if a gene g is ranked i th in both π and σ then $M_{i,i}^F = 2$ and when $\pi = \sigma$, $M^F = 2I$. Let W be a positive semidefinite $m \times m$ matrix, with each entry $W_{i,j}$ being the weight assigned to a gene having rank i and j in the two microarrays. The larger the entries are, the more dependent the two ranks are.

Given these notations we define the distance between the two microarrays as:

$$d(\pi, \sigma) = \sqrt{\sum_{i=1}^m \sum_{j=1}^m ((2I - M_{\pi,\sigma}^F) \circ W)_{i,j}} \quad (13)$$

$$= \sqrt{\sum_{\substack{i,j:\pi^{-1}(i)=\sigma^{-1}(j) \\ \text{or } \pi^{-1}(j)=\sigma^{-1}(i)}} \left(\frac{W_{i,i} + W_{j,j}}{2} - W_{i,j} \right)} \quad (14)$$

$$d(\pi, \sigma) = \sqrt{\text{tr}((2I - M_{\pi,\sigma}^F) W)} \quad (15)$$

where \circ is the Hadamard, or simply entry-wise product. As mentioned above, if the two permutations are identical then $M^F = 2I$ and the distance is 0. Otherwise, the penalty for a disagreement of a pair (i, j) between the rankings is $(W_{i,i} + W_{j,j}) / 2 - W_{i,j}$. This captures both the importance of the individual ranks (very high or very low ranking genes maybe more important than middle genes) as well as the penalty for the disagreement between the pair. Equation (14) also shows that the entity under the square root is non-negative since for a positive semidefinite matrix W , $(W_{i,i} + W_{j,j}) / 2 \geq W_{i,j}, \forall i, j$. Equation (15) follows from Equation (13) since M^F has only one entry in each column / row. In Supporting Methods we prove that this distance function is a pseudometric in the original permutation space and a metric in the W -transformed space.

Learning algorithm: To determine the values of W using the training data we solve the following opti-

mization problem:

$$\min \sum_{(x,y) \in \mathcal{S}} \text{tr} \left(\left(2I - M_{\pi_x, \pi_y}^F \right) W \right) \quad (16)$$

$$\text{subject to } \sum_{(x,y) \in \mathcal{D}} \text{tr} \left(\left(2I - M_{\pi_x, \pi_y}^F \right) W \right) = 1 \quad (17)$$

$$W \succeq 0 \quad (18)$$

Like for the weight vector the constraint (equality to 1) is arbitrary and guarantees that dissimilar arrays are distant from each other. This optimization is a semidefinite program (SDP) [20]. The objective function is a summation of traces of semi-definite matrices and so this is a convex optimization problem and there exists a global minimum solution. However, the matrix W is very large (m by m) and would require large amounts of training data for learning. Since such data is limited using a full rank matrix will likely lead to overfitting. Instead we seek a low-rank approximation of W . Let Z be the rank k approximation of W : $W \approx Z = YY^T$, where $Y \in R^{n \times k}$. Given these changes the optimization problem is:

$$\min \text{tr} (Y^T Z_{\mathcal{S}} Y) \quad (19)$$

$$\text{subject to } \text{tr} (Y^T Z_{\mathcal{D}} Y) = 1 \quad (20)$$

with

$$Z_{\mathcal{S}} = \sum_{(x,y) \in \mathcal{S}} (M_{\pi_x} - M_{\pi_y})^T (M_{\pi_x} - M_{\pi_y})$$

$$Z_{\mathcal{D}} = \sum_{(x,y) \in \mathcal{D}} (M_{\pi_x} - M_{\pi_y})^T (M_{\pi_x} - M_{\pi_y})$$

See Supporting Methods for a discussion on how to further regularize this optimization problem and how to solve it using augmented Lagrangian approach.

3 Experiments and Results

We first used a training dataset from human and mouse tissues to learn parameters for our distance functions and to test the different methods on a dataset for which the correct answer is known. We next downloaded a large number of microarray expression datasets from GEO and applied our distance function to select pairs of experiments that are similar. For this section we consider the cross-species analysis between human (*Homo sapiens*) and mouse (*Mus musculus*) biological samples. We obtained the list of 16,376 human and mouse orthologs from Inparanoid (inparanoid.sbc.su.se).

3.1 Gene Variance

While the methods described above can work for any number of orthologs, the larger the number the more data we would need to fit the weight vector and matrix methods. Since all our expression levels were log ratios to a reference data (see below) we have excluded from the analysis genes that did not vary much *within* each species. We selected the top 500 most varying orthologs for further analysis. We note two things. First, methods that are not affected by overfitting (in our case Spearman's correlation and TBR) were also tested using all orthologs with results very similar to the results obtained from the 500 gene list. Second, while such a selection favors genes with high variance across a large number of experiments, at no stage in the selection have we considered the agreement between the actual levels of orthologous genes in specific experiments.

3.2 Testing distance metrics on data from human and mouse tissues

For evaluation and comparisons of all metrics discussed in this paper, we used an expression dataset, which we call ‘Toronto dataset’, consisting of expression profiles for 26 human tissues and their corresponding tissues in mice [6]. These 26 tissues pairs were profiled using species specific custom arrays. For each tissue, we had one human and one mouse arrays, which were processed and normalized by the authors of [6]. See Supporting Table 1 and Website for the list of tissues.

We used 2 fold cross-validation with 10 random permutations of tissues to compare the performance of the NMF method [25] and the five different distance metrics discussed above. For Pearson correlation, we select the varying 500 genes based on their expression values. For NMF we used the R code provided by the authors which also performs model selection to limit the number of metagenes [4]. The human samples were used to discover the metagenes and the mouse orthologs of these genes were used for the mouse metagenes. For training of the methods discussed in this paper we use the set of similar tissues as the positive set and all the remaining pairs as negative examples. Using parameters learned in the training phase we rank all test pairs by their distance and plot a Precision-Recall (PR) curve for all methods. Since the data set is highly skewed (i.e. there are many more negative than positive pairs), PR curves provide a more informative picture of the metrics’ performance than the Receiver Operator Characteristic (ROC) curves [7].

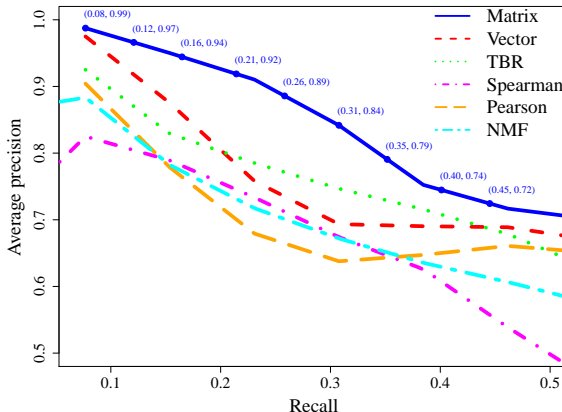


Figure 2: Comparison of different metrics using human-mouse tissues. PR curves of Spearman’s rank correlation, TBR, NMF, Vector and Matrix Weight metrics.

3.2.1 Comparison of cross species comparison metrics

As can be seen in Fig. 2 most methods (except for Spearman’s rank correlation) achieved a very high precision to begin with (80% and higher). However, this precision level drops and when reaching 20% recall only the weight matrix method achieves a precision that is higher than 90%. Since there are hundreds of thousands of expression experiments in GEO, precision is more important than recall for our goals. At these high precision rates the weight matrix method dominates the other methods we have considered and thus we used it in all subsequent analysis.

As for the other methods we believe that Spearman’s rank correlation performs worse than Pearson correlation because the test dataset is well normalized so nonparametric methods loose statistical power. However, in application to large, heterogenous, datasets the assumption of normalization across the datasets is less likely. For NMF, the fact that it is unsupervised and does not use information from the query species to construct the components likely led to its weaker performance. The results presented in Fig. 2 used an approximation matrix with rank 3. We have also tested other ranks (recall that rank 1 is the weight

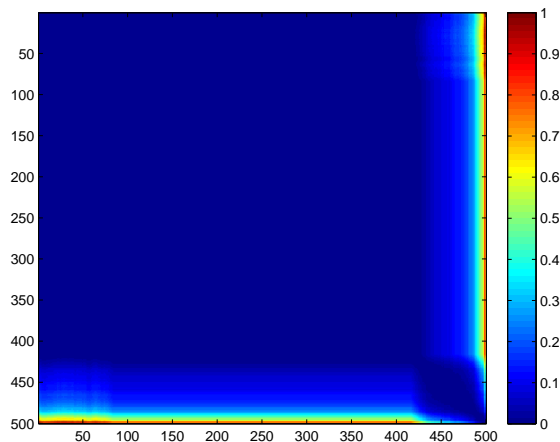


Figure 3: The penalty matrix between ranks $(W_{i,i} + W_{j,j})/2 - W_{i,j}$ as shown in (14), learned from the human and mouse tissues data.

vector shown on the figure as well). We observe that both ranks 2 and 4 do not improve the overall success (Website) and so we have focused on rank 3 matrices for the remainder of this paper.

We have repeated the above analysis (comparison of methods) using another, independent, human-mouse tissue dataset, which we term the ‘Novartis dataset’, from [24]. As we discuss in Supporting Results this additional analysis agrees with the results presented above indicating that our method is robust to the specific data used and to the different platforms in these two studies.

Fig. 3 presents the residual weights $(W_{i,i} + W_{j,j})/2 - W_{i,j}$ which are the penalties for differences in a ranked pair as shown in (14). High (red) values indicate bigger penalty while lower (blue) values indicate that the penalty is smaller. Interestingly the method seems to focus more on the repressed genes and puts a higher weight on genes that move from being repressed to being upregulated or at a medium expression level.

3.2.2 Effect of ortholog assignment on the performance of the Matrix method

Inparanoid contains over 10000 known orthologs between human and mouse making them one of the best annotated pairs of species. As noted above, from this set we select a subset of 500 genes and use these in our algorithms. To test whether our methods would be appropriate to other species pairs for which much fewer orthologs are known we repeated the analysis discussed above starting with a smaller set of orthologs. We selected random sets of 2000 orthologs (roughly 12% of all orthologs) and then reran our method using this initial set (selecting the top 500 varying genes from this smaller subset and running the matrix algorithm discussed above). Figure 4 presents results for seven of these random sets. The blue curve are the results when starting with the full set of orthologs. As can be seen our method is robust and is appropriate for pairs of species with much fewer known orthologs as well.

3.3 Identifying similar experiments in GEO

The previous section shows that our weight matrix performs better than standard metrics on the Toronto and Novartis datasets and moreover can get a very high precision for the recall value of 20%. Our goal is to apply this new metric for retrieving cross-species similar pairs of microarray experiments in a large dataset.

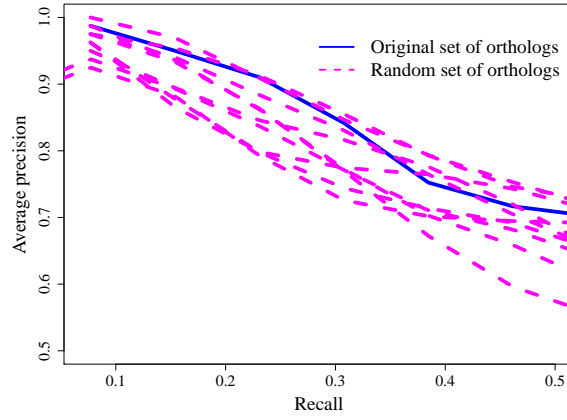


Figure 4: PR curves for the Matrix Weight metric when starting with fewer orthologs. The blue curve is the result when starting with all orthologs (same curve as in Fig. 2).

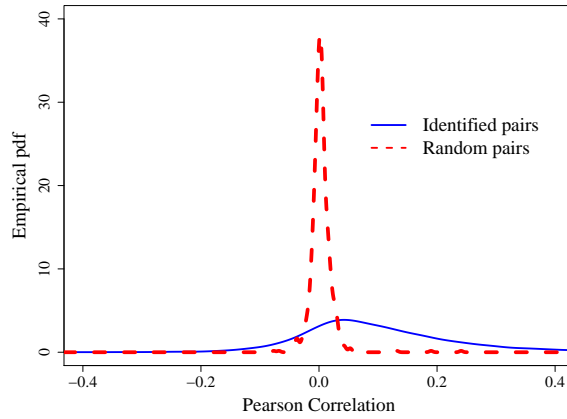


Figure 5: Blue curve: Correlation of orthologs not used for training in a random sample of 301,453 microarray pairs from human and mouse. Red curve: Correlation of orthologs not used for training in the set of microarray pairs selected by our method.

3.3.1 Data Collection

We downloaded 715 human and 769 mouse datasets from GEO and used GDS data and metadata to identify control samples for each dataset (Website). Such samples are important for properly normalizing and transforming the data so that all data used is log2 ratio of the response sample to its control. We excluded from the analysis all datasets for which we could not positively identify the control sample leaving us with 3416 human and 2991 mouse microarrays from 535 human and 641 mouse datasets.

3.3.2 Identification of associated pairs of microarrays

We used the weight matrix trained using the full set of human-mouse tissue pairs. We used the results of Fig. 2 to select a similarity cutoff corresponding to the cutoff that led to 95% precision and 10% recall. Using this cutoff we ended up with 301,453 pairs of microarrays whose distances are smaller than the cutoff which is roughly 3% of all pairs tested. These pairs are from 14493 dataset pairs (many array pairs are from the same pair of human and mouse datasets).

We also looked at the distribution of scores under the null hypothesis (since more than 95% of microarray pairs are not similar, this can be done by selecting random human-mouse array pairs) and determined that the p-value for the null hypothesis is uniformly distributed, as expected. As a sanity check for our results we also computed the Pearson correlation across the pairs determined to be significant by our method for all human and mouse orthologs that were not part of the 500 genes we used for learning the parameters. Fig. 5 shows the histogram of this correlation and the histogram of the correlation for the same set of genes in a randomly selected set of 301,453 microarray pairs. As can be seen the selected experiments are indeed more similar for many of the orthologs when compared to random selected pairs indicating that our method can identify correlated array pairs without using the experiment description.

3.3.3 Description and dataset analysis

The list of pairs derived by our method allows us to address many questions. We first asked what conditions / organs / tissues are the most similar between human and mouse in terms of expression. We used the titles provided in the metadata section of the GDS to identify common words that are significantly over-represented in the microarray pairs we extracted. For each pair of similar experiments, a word that appears in both titles could provide information about the relationship between the pair. For each word we have also computed the number of times it appeared in a title for all microarrays used from each species and the expected number of times it should have appeared in the pairs we selected. Using the hypergeometric distribution we computed the overrepresentation P-value for each word. Table 1 presents the results of the analysis of over represented words in matched titles. As can be seen some organs and tissue types are much more represented than others. For example, brain, muscles and blood appear to have similar expression patterns between the two species. Certain conditions are also overrepresented, most notably immune response. Several words are associated with experiments related to such response including different types of cells participating in the response (macrophages, dendritic, cd8). In contrast, cancer, one of the most common words in the human studies (roughly 10% of human datasets contained cancer in the title) was not overrepresented supporting recent results that most mice are not an ideal model system for at least some types of cancer [3, 22]. We repeated this analysis using the abstracts provided instead of the titles leading to similar results (see Website for full results). We have also looked beyond pairwise similarities and identified entire datasets (GDS files) that contained several similar pairs of arrays between human and mouse. An expert pathologist (Oltvai) manually inspected the top 100 matched datasets and determined that over 80% of them make biological sense (see Supporting Table 2). Many of the datasets identified as similar contained experiments for the same tissue (most notably muscle, but also blood and brain). However, some of the matches were less obvious. Fibrosis is a chronic progressive and often lethal lung disease. One of the top 50 matches in our results was between a human dataset titled non-diseased lung tissue (GDS1673) and the mouse dataset titled Pulmonary fibrosis(GDS251). However, upon a closer inspection of the mouse dataset it can be seen that it compares two mouse strains treated with bleomycin. One is determined to be susceptible to fibrosis (C57BL6/J) whereas the other is determined to be resistant (BALB/c). When looking at the similarities computed by

Rank	P-value	Word	#Pairs	
			Identified	Expected
1	7.14429e-13	MUSCLE	121	28.46752
2	7.39409e-13	DENDRITIC	24	2.13506
3	1.76946e-11	SKELETAL	42	12.12506
4	3.12418e-11	MACROPHAGE	18	2.21414
5	1.89634e-08	ERYTHROID	6	0.15815
6	2.52933e-08	OBESITY	9	0.63261
7	8.35063e-08	HEMATOPOIETIC	13	1.84512
8	2.36749e-07	BRAIN	19	4.42828
9	1.52768e-06	CD8+	5	0.18451
10	1.67619e-06	CARDIAC	6	0.34266
11	1.45374e-05	STEM	43	20.87618
12	2.02795e-05	HAIR	5	0.31631
13	9.19217e-05	FIBROBLASTS	12	3.08398
14	2.04560e-04	AIRWAY	7	1.15979

Table 1: Top 14 words identified in titles of pairs determined to be similar. #Pairs Identified is the number of time this pair was observed. #Pairs Expected is the number of time expected based on single species occurrences. The P-value is computed using the hypergeometric distribution.

Rank	Category Name	# Genes			
		Assigned	Expected	P	P adj
1	cell cycle	39.0	9.1	8.5E-15	<0.001
2	cell division	26.0	4.5	5.5E-13	<0.001
3	cell cycle phase	26.0	4.7	1.6E-12	<0.001
4	M phase	24.0	4.2	4.8E-12	<0.001
5	cell cycle process	26.0	5.5	4.6E-11	<0.001
6	mitotic cell cycle	21.0	3.8	2.4E-10	<0.001
7	mitosis	17.0	2.9	6.7E-9	<0.001
8	nuclear division	17.0	2.9	5.8E-9	<0.001
9	M phase of mitotic cycle	17.0	3.0	6.7E-9	<0.001

Table 2: GO enrichment analysis for mouse genes using STEM.

our method it can be seen that the vast majority of the top 100 matches are for the BALB/c strains. Thus, our cross species comparisons can be used to identify cases in which similar pathways are activated even though the conditions may be different.

3.3.4 Quarrying GEO to identify cycling mouse genes

To demonstrate the utility of our method for quarrying large cross species databases like GEO we used a set of 50 known human cycling genes extracted from [26]. For each of these genes we used all 301,453 microarray pairs determined to be similar to identify the set of similarly expressed mouse genes using Spearman correlations (regardless of their sequence similarity). We retrieved the top 10 most similar mouse genes for each query human gene resulting in a set of 206 genes. Note that the database we used contained a diverse set of experiments and, while a few may have been focused on cell cycle studies the vast majority were not. Importantly, our analysis here did not rely on any specific cell cycle time series dataset.

We used STEM [9] to determine significant GO categories associated with this list of mouse genes. As can be seen in Table 2, all top categories that are enriched for this set are related to cell cycle (including cell cycle itself). The set of mouse genes contains orthologs of the original set of human genes including CDC2A, a cell division control protein and CCNB1, an essential component of the cell cycle regulatory machinery. The list also contains many known mouse cell cycle genes with no homologs on the human list. These include members of a highly conserved complex which is essential for the initiation of DNA replication (ORC1L and ORC6L) and PRIM1 and PRIM2 which are involved in chromosomal replication during cell

cycle. See Website for complete list. These results highlight the potential use of our method for identifying functionally related genes across species.

4 Conclusions and future work

The growth of microarray databases opens the door to applications that can simultaneously query sequence and expression databases to identify both static and dynamic matches. However, these methods would require a set of matching expression datasets in the species being queried. Such matches are hard to come by. It is rare to find the exact same experiment (condition, time, tissues etc.) in multiple species. To allow the use of these databases we looked at several different distance metrics between expression experiments. We defined a new distance function which utilizes the ranking of orthologs in both species. Our method uses a training dataset to learn weights for differences in rankings between the species and these differences are then summed up to determine the similarity between the two experiments. Testing this method on a training dataset of known similar pairs showed that it indeed improves upon other distance measures and that it can achieve high precision.

We used our new distance function to retrieve similar experiment pairs from GEO. The set of experiments identified by our method allowed us to look at questions regarding the conditions and tissues that activate similar expression patterns in human and mouse and to find a set of cycling mouse genes based on a set of known human cycling genes. Many of these mouse genes are known to be cycling and the rest of the genes identified are candidates for further study into their role in the cell cycle.

Our method attempts to learn a new distance function for permutations based on training data. There has been recent work in Machine Learning on trying to learn new distance function for feature vectors [2], though we are not aware of any work so far that attempted to learn such methods for permutations. A number of the methods developed for feature vectors were later kernelized allowing for much faster computations. It would be interesting to see if the Matrix weight method discussed in this paper can also be kernelized. We have primarily relied on one to one orthology matches for computing the distance between pairs of experiments. Since many orthology assignments are many to one or many to many, methods that can utilize such information may be able to improve upon the results suggested in this paper. Our overall goal is to compile a large set of expression pairs that can be used for querying human and mouse genes. As we noted in the introduction our method can also help in distinguishing between orthologs and homologs by looking for genes with similar sequence that are also co-expressed in the set of similar experiments. We would also like to extend this work to other species and we are looking for training data for these species.

Appendix

A Metric properties

For completeness we list below the properties of distance metrics.

1. Non-negative: $d(\pi, \sigma) \geq 0$
2. Symmetric: $d(\pi, \sigma) = d(\sigma, \pi)$
3. Identity: $d(\pi, \sigma) = 0$ if and only if $\pi = \sigma$
4. Triangular inequality: $d(\pi, \sigma) \leq d(\pi, \tau) + d(\tau, \sigma)$ for any $\tau \in G_m$

B Proof of Asymptotic Normality

Proof Since $d(\pi, \sigma) = d(\pi\pi^{-1}, \sigma\pi^{-1}) = d(I_m, \sigma\pi^{-1})$, the distribution D_d is the distribution of $d(I_m, \tau)$ when τ is a uniformly random permutation in G_m . Applying Hoeffding's Combinatorial Central Limit Theorem [12] with $c_m(i, j) = (w(i) - w(j))^2$, we only need to verify the condition (12) of the theorem 3.

Define $d_m(i, j)$ as in the equation (11). Let $\alpha = \min_{1 \leq i, j \leq m} d_m(i, j)$ and $\beta = \max_{1 \leq i, j \leq m} d_m(i, j)$. α and β exist because $-\infty < w(i) < \infty$ for all i .

$$\lim_{m \rightarrow \infty} \frac{\max_{1 \leq i, j \leq m} [d_m(i, j)]^2}{\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^m [d_m(i, j)]^2} \leq \lim_{m \rightarrow \infty} \frac{\beta^2}{\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^m \alpha^2} \quad (21)$$

$$= \lim_{m \rightarrow \infty} \frac{\beta^2}{m\alpha^2} \quad (22)$$

$$= 0 \quad (23)$$

C Pseudometric properties of the relational weighted rank matrix

Below we prove that Equation 13 is a pseudometric in the original permutation space and a metric in the W -transformed space.

Lemma C.1

$$M_{\pi, \sigma}^F = M_{\pi}^T M_{\sigma} + M_{\sigma}^T M_{\pi} \quad (24)$$

$$2I - M_{\pi, \sigma}^F = (M_{\pi} - M_{\sigma})^T (M_{\pi} - M_{\sigma}) \quad (25)$$

Proof Since M_{π} and M_{σ} are permutation matrices, $M_{\pi}^T M_{\sigma} = M_{\sigma\pi^{-1}}$ and $M_{\sigma}^T M_{\pi} = M_{\pi\sigma^{-1}}$.

Therefore, by the definition of the permutation matrix in (5), $M_{\sigma\pi^{-1}}(i, j) = 1$ if and only if $\sigma\pi^{-1}(i) = j$ or $\pi^{-1}(i) = \sigma^{-1}(j)$. Similarly, $M_{\pi\sigma^{-1}}(i, j) = 1$ if and only if $\pi^{-1}(j) = \sigma^{-1}(i)$. Equation (24) follows from the definition of $M_{\pi\sigma}^F$ in (12).

$$\begin{aligned} 2I - M_{\pi, \sigma}^F &= M_{\pi}^T M_{\pi} + M_{\sigma}^T M_{\sigma} - (M_{\pi}^T M_{\sigma} + M_{\sigma}^T M_{\pi}) \\ &= (M_{\pi} - M_{\sigma})^T (M_{\pi} - M_{\sigma}) \end{aligned}$$

Theorem C.2 *If the matrix W is positive semidefinite, the distance is a pseudometric.*

Proof

$$\begin{aligned} d(\pi, \sigma) &= \sqrt{\text{tr}((2I - M_{\pi, \sigma}^F)W)} \\ &= \sqrt{\text{tr}(Y^T(M_\pi - M_\sigma)^T(M_\pi - M_\sigma)Y)} \\ &= \|(M_\pi - M_\sigma)Y\|_F \end{aligned}$$

Since the Frobenius norm $\|\cdot\|_F$ is a metric, our distance $d(\pi, \sigma)$ satisfies non negativity, symmetry and triangular inequality. Therefore, the distance is a pseudometric. $d(\pi, \sigma) = 0$ implies $M_\pi Y = M_\sigma Y$, hence the distance is a metric in the W -transformed space.

D Matrix and Vector Weight metrics

We show that the vector weight discussed in section 2.3.3 is a special case of the general weight matrix when that matrix has a rank of 1.

Proof Since W is ranked 1, $W = w^T w$ with w is a vector of length n . Let d_1 and d_2 be the metric in Sect. 2.3.3 and Sect. 2.3.4 respectively. Recall from the proof of Theorem C:

$$d_1(\pi, \sigma) = \sqrt{w^T(M_\pi - M_\sigma)^T(M_\pi - M_\sigma)w} \quad (26)$$

$$d_2(\pi, \sigma) = \sqrt{\text{tr}(w^T(M_\pi - M_\sigma)^T(M_\pi - M_\sigma)w)} \quad (27)$$

$$= \sqrt{w^T(M_\pi - M_\sigma)^T(M_\pi - M_\sigma)w} \quad (28)$$

Therefore, the metric in Sect. 2.3.3 is a special case of the metric in Sect. 2.3.4.

E Regularizing and solving the weight matrix optimization problem

An additional constraint that is useful for controlling overfitting is to regularize the solution. In our case, since nearby locations can be affected by small amounts of noise a reasonable regularization policy is to require that the W matrix is smooth. To achieve this we add linear inequality constraints to enforce that column-adjacent entries in Y differ by at most $\delta > 0$: $Y(i, j) - Y(i+1, j) \geq -\delta$ and $-Y(i, j) + Y(i+1, j) \geq -\delta$ $\forall 1 \leq i < m, 1 \leq j \leq k$.

We solve this optimization by using the augmented Lagrangian approach. Similarly, we can incorporate the smoothness constraints to the Lagrangian. See [20] for a detailed discussion on the augmented Lagrangian method.

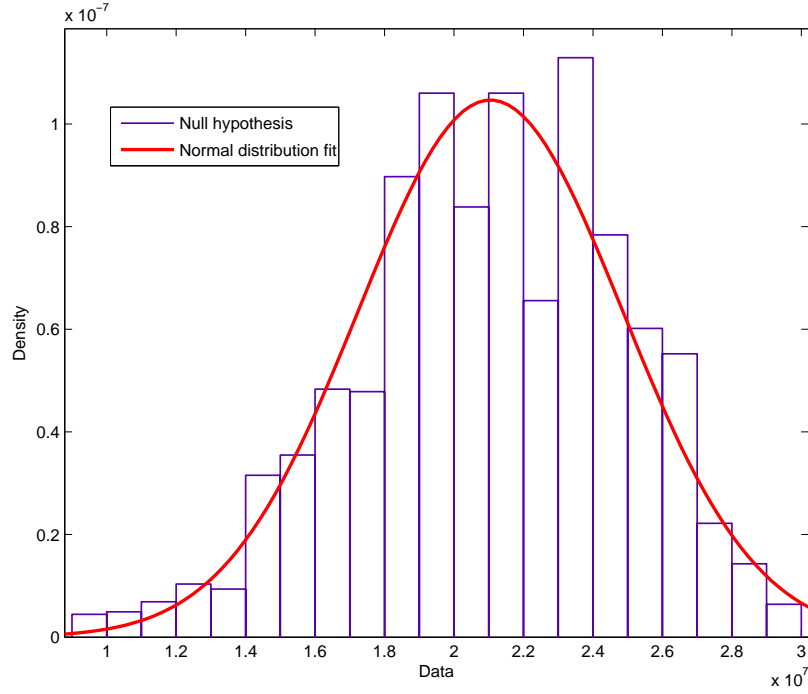
F Normality of the null distribution

Supporting Fig. 1 experimentally confirms that the null model follows a normal distribution. The red curve is a normal distribution fit using Matlab.

G Testing distance metrics data from human and mouse tissues

G.1 Different rank values and number of negative examples

Supporting Fig. 2 shows the PR curves of the Matrix Weight metrics using the rank values of 2,3 and 4. Both ranks 2 and 4 do not improve the overall success. We also have tested using a different number of



Supporting Fig. 1: The histogram of the Spearman correlation of 2000 random pairs of microarrays and the Gaussian distribution fit using Matlab.

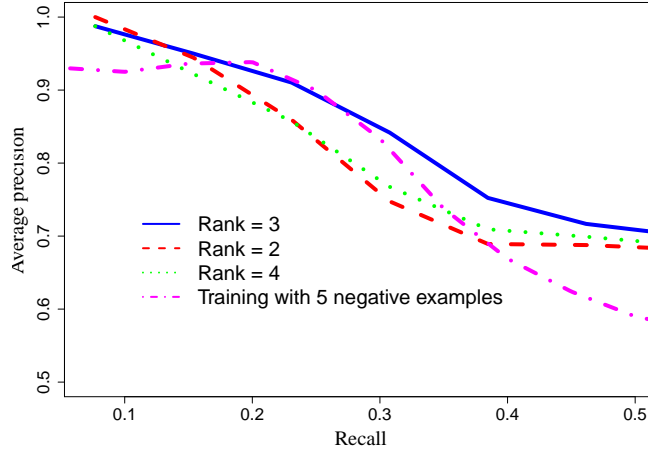
negative examples for each array in the training set (since the number of positive examples is only 1 it is hard to change that number). For this test we used 5 negative examples (in the original analysis we used 12). As can be seen in Supporting Fig. 2, this change did not affect the results much and the PR curve for such setting is very close to the original PR curve.

G.2 Comparison of cross species comparison metrics using 1000 most variant genes

We reran experiments with 1000 orthologs and the results are presented in Supporting Fig. 3 . Indeed, as the reviewer suspected the matrix method did slightly worse when compared to the results using 500 genes. However, for the highest precision rates Matrix was still the best method (though by a much lower margin when compared to the vector method which requires far fewer parameters). The results of using 500 genes are slightly better than using 1000 genes at the 0.9 precision range (for a recall of 0.21 the 500 genes method achieves a 0.92 precision whereas the 1000 genes achieves 0.91). Of course, these results are also a function of the training data size. With a larger training datasets the ability to fit parameters to more sophisticated models increases and so more complex methods, like the Matrix method, are likely to outperform the simpler methods.

G.3 Randomized dataset

To demonstrate that how well different methods perform relative to random prediction, we have carried out the experiment on a randomized dataset, by randomly permuting expression values in each array. The results are presented as Supporting Fig. 4. As can be seen, all methods do very badly and the results are essentially a flat PR curve as expected from random data.



Supporting Fig. 2: PR curves of Matrix Weight metrics with different rank values and using less number of negative examples.

H Testing distance metrics on an additional dataset from human and mouse tissues

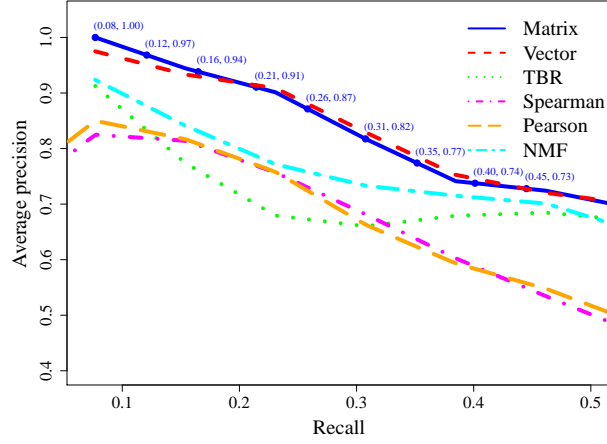
For an additional evaluation of all metrics discussed in this paper, we used a second human-mouse expression dataset consisting of 79 human and 61 mouse tissues from [24] (note that some are repeats). In cases where the cell types differed between human and mouse we have assigned each human tissue sample to at most three mouse samples based on a mapping by a pathologist (Oltvai). The assignment of human tissues to mouse tissues are based on the following criteria (see Website for complete assignments):

1. Same organs, cell types, and developmental stages.
2. Spatially closer structures within an organ.
3. Insights that are not necessarily evident from anatomy, e.g, the ontogenic similarity of brown adipose tissue and muscle.

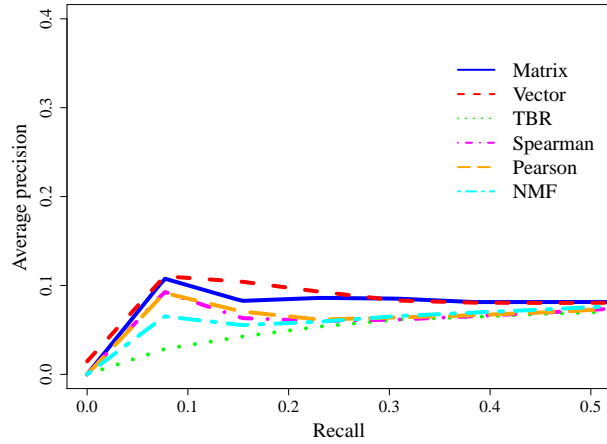
We next used 4 fold cross-validation with 4 random permutations of the tissues to compare the performance of the NMF method [25] and the four different distance matrices discussed above. The results presented used an approximation matrix with rank 3.

The overall success for this dataset is lower than for the Toronto dataset. This agrees with the initial analysis of this data that indicated a large deviation between human and mouse expression data for some of the tissues [24]. Still, in terms of comparison between methods the results of this analysis agree with the results presented in the main text. As can be seen in Supporting Fig. 5(a) the weight matrix method achieves a high precision (65%) for a much larger recall (10%). As discussed in the main text the reason NMF does not perform well on this dataset is likely related to the fact that it is unsupervised and does not use information from the query species to construct the components.

Supporting Fig. 5(b) presents the residual weights $(W_{i,i} + W_{j,j})/2 - W_{i,j}$ which are the penalties for differences in a ranked pair as shown in (14). We note the similarity with the learnt matrix in Fig. 3 in putting a higher weight on genes that move from being repressed although the penalty is smaller. Thus, the overall weighting seems to be dataset and platform independent.



Supporting Fig. 3: PR curves of Spearman's rank correlation, TBR, NMF, Vector and Matrix Weight metrics using 1000 most variant genes.



Supporting Fig. 4: PR curves of Spearman's rank correlation, TBR, NMF, Vector and Matrix Weight metrics on a randomized dataset.

I Human and mouse tissue list

Supporting Table 3 shows the list of 26 human and mouse tissues used in this analysis.

J Identifying similar experiments in GEO

J.1 Histogram of the correlation of 500 selected genes

Supporting Fig. 6 shows distributions of correlations for the selected highly varying 500 genes. When using the 500 selected genes the results look pretty similar to the results presented in the paper though the mean correlation is slightly higher (0.1057 vs. 0.1021).

Human	Mouse
Adrenal Cortex	Adrenal
Bladder	Bladder
Bone Marrow	Bone Marrow
Brain	Brain
Brain Cerebellum	Cerebellum
Brain Cerebral cortex	Cortex
Epididymis	Epididymus
Heart	Heart
Kidney	Kidney
Liver	Liver
Lung	Lung
Pancreas	Pancreas
Placenta	Placenta 12.5
Prostate	Prostate
Salivary Gland	Salivary
Skeletal Muscle	Skeletal Muscle
Small Intestine	Small Intestine
Spinal Cord	Spinal Cord
Spleen	Spleen
Stomach	Stomach
Testis	Testis
Thymus	Thymus
Thyroid	Thyroid
Tongue	Tongue
Trachea	Trachea
Uterus	Uterus

Supporting Table 3: The one-one similarity list of human and mouse tissues.

J.2 Description analysis on random sets of array pairs

We repeated the analysis with random sets of array pairs. As can be seen in Supporting Table 4, for these pairs the p-values are much higher (less significant). Specifically, there are no matched terms with a p-value lower than 10^{-10} (whereas in the identified matching there are 4 such words) and only 3 of the top random match words would be ranked in the top 10 of the words identified using the matches made by the algorithm. Thus, such p-values are significant and would not be expected from random assignments.

J.3 Heat map of similarity between 3416 human and 2991 mouse microarrays

Supporting Fig. 7 presents a heatmap showing all human by mouse arrays where the color indicates the level of similarity from the Weight Matrix metric. Smaller value means more similarity.

J.4 Human assessment of identified matched dataset pairs.

To test whether the identified matched pairs are indeed a feasible solution we have asked an expert pathologist (Oltvai, a co-author of the paper) to examine the top 100 matched dataset pairs identified by our method. Based on the description for that dataset the expert assigned each match to one of three categories: A correct match (Y), an incorrect match (N) and an inconclusive. As can be seen in Supplementary Table 5, there were 83 Y assignments in the top 100 matches with the other 17 determined to either be mistakes (N, 13) or inconclusive (4). Given that almost all random matches would not make sense this is a very high accuracy

Rank	P-value	Word	#Pairs	
			Identified	Expected
1	1.13469e-09	BONE	51	19.13650
2	3.91648e-09	ACUTE	43	15.18268
3	1.26953e-06	MARROW	15	3.16306
4	1.49012e-05	GASTROCNEMIUS	8	1.05435
5	2.02795e-05	STEROID	5	0.31631
6	7.76604e-05	METAPLASIA	3	0.07908
7	1.34712e-04	LIPOPOLYSACCHARIDE	8	0 1.44973
8	2.29396e-04	PULMONARY	15	05.00818
9	3.32228e-04	PROGENITOR	8	0 1.66061
10	5.00167e-04	IFN-GAMMA	5	0.63261
11	7.80427e-04	DYSTROPHY	14	5.06089
12	7.86850e-04	DUCHENNE	8	1.89783
13	7.94474e-04	REGIONS	9	02.37229
14	1.34160e-03	LEUKEMIAS	2	0.05272

Supporting Table 4: Top 14 words identified in titles of pairs determined to be similar. #Pairs Identified is the number of time this pair was observed. #Pairs Expected is the number of time expected based on single species occurrences.

rate and it clearly indicates that this method can be use to help improve, and speed up, human assessment of similarity. We have changed the introduction and results sections to reflect this idea and to highlight the ability of the method to aid in human assessment of similarity.

Human Dataset	Description	Mouse Dataset	Description	Assessment
GDS2767	Blood response to various beverages: time course	GDS1077	Hematopoietic stem cells from different recombinant inbred strains	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS2047	Lipopolysaccharide effect on macrophages pretreated with carbon monoxide: time course	Y/ inconcl.
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS2330	Acute myocardial infarction model: time course (MG-U74B)	Y
GDS2055	Skeletal muscle types (HG-U133A)	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y
GDS1815	High-grade gliomas (HG-U133A)	GDS2159	Spinal cord injury model: time course	Y
GDS2055	Skeletal muscle types (HG-U133A)	GDS2330	Acute myocardial infarction model: time course (MG-U74B)	Y
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS488	Myocardial infarction time course	Y
GDS2056	Skeletal muscle types (HG-U133B)	GDS2330	Acute myocardial infarction model: time course (MG-U74B)	Y
GDS2740	Lengthening and shortening contractions effect on the muscle: time course	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y
GDS2740	Lengthening and shortening contractions effect on the muscle: time course	GDS2330	Acute myocardial infarction model: time course (MG-U74B)	Y

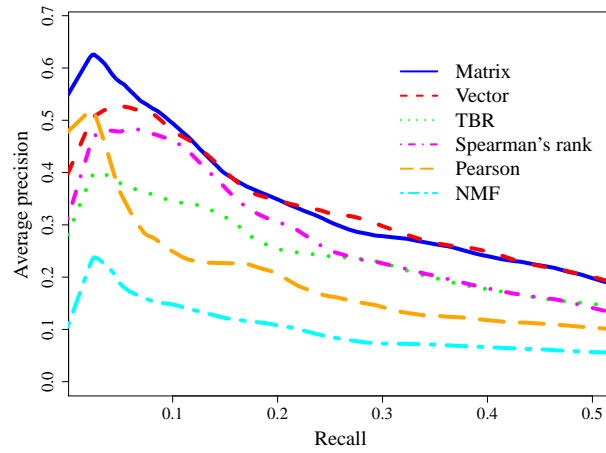
GDS2678	Brain regions of humans and chimpanzees	GDS2159	Spinal cord injury model: time course	Y
GDS2055	Skeletal muscle types (HG-U133A)	GDS488	Myocardial infarction time course	Y
GDS2767	Blood response to various beverages: time course	GDS2150	Spleens of males and females at puberty	N
GDS2255	Transmigrated neutrophils in the alveolar space of endotoxin-exposed lung	GDS1077	Hematopoietic stem cells from different recombinant inbred strains	Y
GDS2255	Transmigrated neutrophils in the alveolar space of endotoxin-exposed lung	GDS2047	Lipopolysaccharide effect on macrophages pretreated with carbon monoxide: time course	Y
GDS2373	Squamous cell lung carcinomas	GDS2334	Myod and Myog expression effect on myogenesis: time course	N
GDS2373	Squamous cell lung carcinomas	GDS981	Uterine response to physiologic and plant-derived estrogen: time course	N
GDS2373	Squamous cell lung carcinomas	GDS1244	Phosgene effect on lungs: time course	Y
GDS2055	Skeletal muscle types (HG-U133A)	GDS234	Muscle regeneration (U74Av2)	Y
GDS2767	Blood response to various beverages: time course	GDS1336	T cell anergy induction regulation by Egr-2 and Egr-3 (MG-U74A)	Y/ inconcl.
GDS1673	Non-diseased lung tissue	GDS1244	Phosgene effect on lungs: time course	Y
GDS2373	Squamous cell lung carcinomas	GDS1072	Platelet derived growth factor effect in the presence of Src family kinase inhibitors (MOE430A)	Inconcl.
GDS2740	Lengthening and shortening contractions effect on the muscle: time course	GDS488	Myocardial infarction time course	Y
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS2335	Exercise effect on the diabetic cardiac muscle: time course	Y
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS627	Cardiac development in embryo	Y
GDS2767	Blood response to various beverages: time course	GDS1514	Interferon-gamma tolerogenic effect on CD8+ dendritic cells	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS2408	B cell-activating factor of the TNF family effect on B cells	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS993	Naive CD8+ T cells proliferative response to lymphopenia: time course	Y/ inconcl.
GDS2055	Skeletal muscle types (HG-U133A)	GDS1541	Exercise effect on diabetic skeletal muscle: time course	Y
GDS2056	Skeletal muscle types (HG-U133B)	GDS1541	Exercise effect on diabetic skeletal muscle: time course	Y
GDS2678	Brain regions of humans and chimpanzees	GDS2917	Various brain regions of several inbred strains	Y
GDS2055	Skeletal muscle types (HG-U133A)	GDS2335	Exercise effect on the diabetic cardiac muscle: time course	Y
GDS596	Large-scale analysis of the human transcriptome (HG-U133A)	GDS2159	Spinal cord injury model: time course	Inconcl.
GDS2678	Brain regions of humans and chimpanzees	GDS1406	Brain regions of various inbred strains	Y
GDS2373	Squamous cell lung carcinomas	GDS1058	Uterus response to 17beta-estradiol: time course	N
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS1766	Extraocular and hindlimb skeletal muscle cell differentiation: time course (MG-430B)	Y
GDS1340	Exercise effect on aged muscle	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y
GDS1340	Exercise effect on aged muscle	GDS2330	Acute myocardial infarction model: time course (MG-U74B)	Y
GDS198	Inflammatory myopathy	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y

GDS2373	Squamous cell lung carcinomas	GDS1277	Obliterative bronchiolitis and tracheal allograft	Y
GDS1673	Non-diseased lung tissue	GDS251	Pulmonary fibrosis	Y
GDS2767	Blood response to various beverages: time course	GDS882	Neuromedin U effect on type-2 Th cells: time course	Y/ inconcl.
GDS2168	HIV viremia effect on monocytes	GDS1077	Hematopoietic stem cells from different recombinant inbred strains	Y
GDS707	Aging brain: frontal cortex expression profiles at various ages	GDS2159	Spinal cord injury model: time course	Y
GDS2055	Skeletal muscle types (HG-U133A)	GDS1765	Extraocular and hindlimb skeletal muscle cell differentiation: time course (MG-430A)	Y
GDS2373	Squamous cell lung carcinomas	GDS1631	Osteoblast differentiation (MG-U74A)	N
GDS2373	Squamous cell lung carcinomas	GDS1071	Platelet derived growth factor effect in the presence of Src family kinase inhibitors (MG-U74A)	Y
GDS198	Inflammatory myopathy	GDS234	Muscle regeneration (U74Av2)	Y
GDS2767	Blood response to various beverages: time course	GDS1285	Macrophage response to lipopolysaccharide and CstF-64 overexpression	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS1315	Immune response to suppressive vs. stimulatory immunomodulators	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS1654	Dendritic cell subpopulations: spleen (MG-U74A)	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS2741	TCR-alpha/beta CD8-alpha/alpha intestinal intraepithelial lymphocytes	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS2957	Resting and activated natural killer cells	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS658	Thymocyte selection by agonist	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS827	Acute ethanol administration effect on Toll-like receptor 3 signaling in macrophages	Y/ inconcl.
GDS2056	Skeletal muscle types (HG-U133B)	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y
GDS2083	Limb immobilization effect on skeletal muscle	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y
GDS2083	Limb immobilization effect on skeletal muscle	GDS2330	Acute myocardial infarction model: time course (MG-U74B)	Y
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS40	Cardiac development, maturation and aging	Y
GDS2373	Squamous cell lung carcinomas	GDS1865	Chondrocyte differentiation: time course	N
GDS2767	Blood response to various beverages: time course	GDS2521	Megakaryocytes at successive stages of maturation	Y
GDS395	Biomaterial engineering	GDS981	Uterine response to physiologic and plant-derived estrogen: time course	N
GDS2113	Pheochromocytomas of various genetic origins	GDS2159	Spinal cord injury model: time course	Y
GDS1036	Microglial cell response to interferon-gamma: time course	GDS2047	Lipopolysaccharide effect on macrophages pretreated with carbon monoxide: time course	Inconcl.
GDS1684	Cardiac allograft rejection: time course	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Y
GDS1684	Cardiac allograft rejection: time course	GDS2330	Acute myocardial infarction model: time course (MG-U74B)	Y
GDS2740	Lengthening and shortening contractions effect on the muscle: time course	GDS1541	Exercise effect on diabetic skeletal muscle: time course	Y
GDS2740	Lengthening and shortening contractions effect on the muscle: time course	GDS2335	Exercise effect on the diabetic cardiac muscle: time course	Y

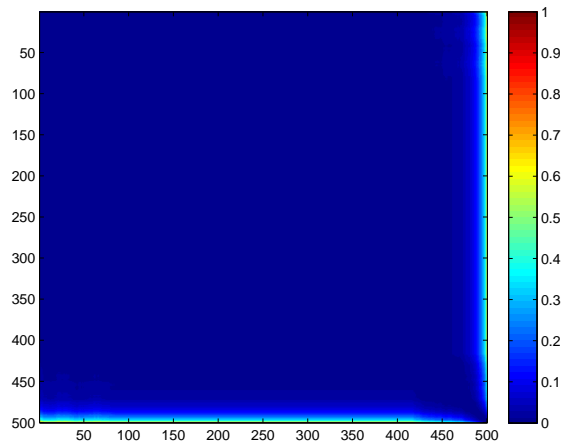
GDS833	Alternative pre-mRNA splicing in various tissues and cell lines (Rosetta/Merck Splicing Chip 5)	GDS2162	CH1 domain deletion, p300 and CBP heterozygous null mutant hypoxic fibroblasts response to trichostatin A	N
GDS833	Alternative pre-mRNA splicing in various tissues and cell lines (Rosetta/Merck Splicing Chip 5)	GDS1244	Phosgene effect on lungs: time course	N
GDS1284	Multiple myeloma molecular classification	GDS1077	Hematopoietic stem cells from different recombinant inbred strains	Y
GDS198	Inflammatory myopathy	GDS488	Myocardial infarction time course	Y
GDS2055	Skeletal muscle types (HG-U133A)	GDS627	Cardiac development in embryo	Y
GDS2373	Squamous cell lung carcinomas	GDS951	Hormone-induced adipogenesis suppressed by 2,3,7,8-tetrachlorodibenzo-p-dioxin and EGF	N
GDS424	Normal human tissue expression profiling (HG-U95C)	GDS2329	Acute myocardial infarction model: time course (MG-U74A)	Inconcl.
GDS2255	Transmigrated neutrophils in the alveolar space of endotoxin-exposed lung	GDS1336	T cell anergy induction regulation by Egr-2 and Egr-3 (MG-U74A)	Y/ inconcl.
GDS2528	Basal plate of the placenta from midgestation to term (HG-U133A)	GDS981	Uterine response to physiologic and plant-derived estrogen: time course	Y
GDS1340	Exercise effect on aged muscle	GDS488	Myocardial infarction time course	Y
GDS2056	Skeletal muscle types (HG-U133B)	GDS488	Myocardial infarction time course	Y
GDS1340	Exercise effect on aged muscle	GDS234	Muscle regeneration (U74Av2)	Y
GDS2373	Squamous cell lung carcinomas	GDS857	Corneal stromal cell differentiation	N
GDS1815	High-grade gliomas (HG-U133A)	GDS2917	Various brain regions of several inbred strains	Y
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS1541	Exercise effect on diabetic skeletal muscle: time course	Y
GDS2106	Lymphoblastoid cell lines from various CEPH pedigrees	GDS2047	Lipopolysaccharide effect on macrophages pretreated with carbon monoxide: time course	Y
GDS2310	Exercise effect on white blood cells	GDS2047	Lipopolysaccharide effect on macrophages pretreated with carbon monoxide: time course	Y
GDS2772	Sevoflurane and propofol effect on the heart during off-pump coronary artery bypass graft surgery	GDS388	Cardiac remodeling (Mu11K-B)	Y
GDS1962	Glioma-derived stem cell factor effect on angiogenesis in the brain	GDS2159	Spinal cord injury model: time course	Y
GDS2255	Transmigrated neutrophils in the alveolar space of endotoxin-exposed lung	GDS1514	Interferon-gamma tolerogenic effect on CD8+ dendritic cells	Y
GDS2255	Transmigrated neutrophils in the alveolar space of endotoxin-exposed lung	GDS2408	B cell-activating factor of the TNF family effect on B cells	Y
GDS738	Intervertebral disc cells and osmotic loading	GDS981	Uterine response to physiologic and plant-derived estrogen: time course	N
GDS395	Biomaterial engineering	GDS2162	CH1 domain deletion, p300 and CBP heterozygous null mutant hypoxic fibroblasts response to trichostatin A	N
GDS2435	Male and female venous blood	GDS1077	Hematopoietic stem cells from different recombinant inbred strains	Y
GDS2959	Granulocyte colony-stimulating factor mobilized leukocytes	GDS1077	Hematopoietic stem cells from different recombinant inbred strains	Y
GDS2767	Blood response to various beverages: time course	GDS2011	Lupus-prone BWF1 males and females: spleen (MG-U74A)	Y/ inconcl.

GDS2767	Blood response to various beverages: time course	GDS2041	Type II activated macrophage	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS2651	Macrophage cell line response to Chlamydia pneumoniae infection	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS433	CD8+ effector and central memory T cells (MG-U74A)	Y/ inconcl.
GDS2767	Blood response to various beverages: time course	GDS684	T regulatory and T effector cells in prediabetic lesion	Y/ inconcl.
GDS2055	Skeletal muscle types (HG-U133A)	GDS2001	Utrophin/dystrophin-deficient double mutant and dystrophin-deficient mdx mutant skeletal muscles	Y

Supporting Table 5: The result of human assessment of identified matched dataset pairs.

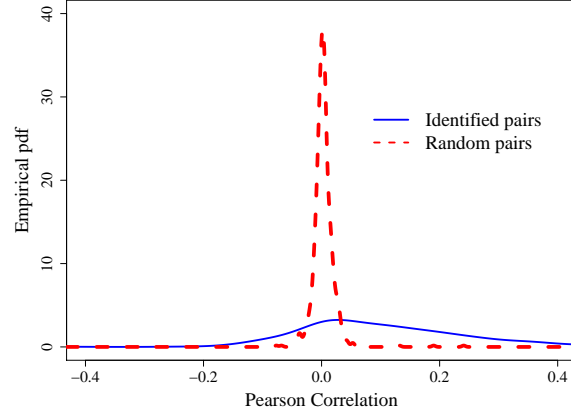


(a) PR curves of Spearman's rank correlation, TBR, NMF, Vector and Matrix Weight metrics.

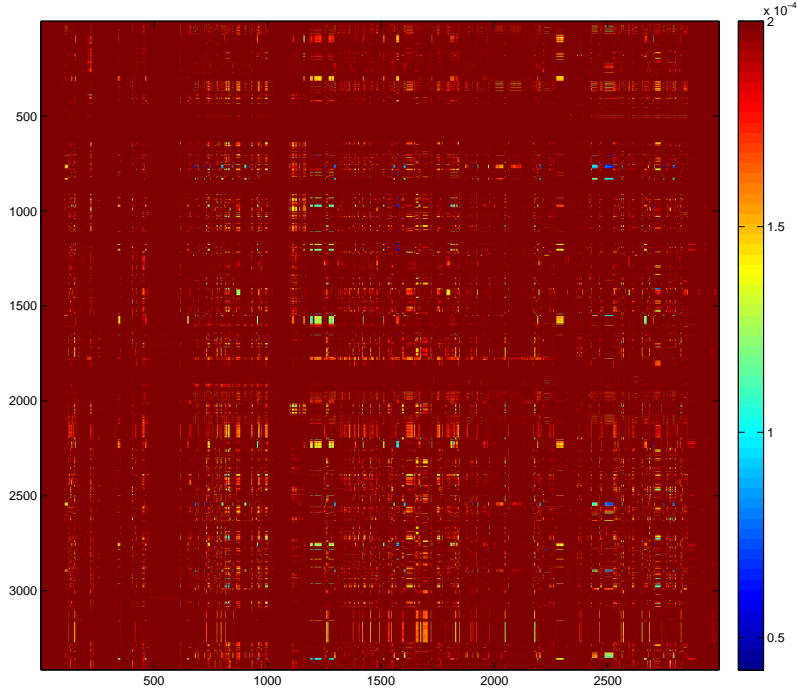


(b) The penalty matrix between ranks $(W_{i,i} + W_{j,j})/2 - W_{i,j}$ as shown in (14), learned from the human and mouse tissues data.

Supporting Fig. 5: Experimental evaluation of metrics.



Supporting Fig. 6: Blue curve: Correlation of 500 orthologs used for training in a random sample of 301,453 microarray pairs from human and mouse. Red curve: Correlation of 500 orthologs used for training in the set of microarray pairs selected by our method.



Supporting Fig. 7: The similarity between 3416 human and 2991 mouse microarrays.

References

- [1] S. F. Altschul, T. L. Madden, A. A. Schäffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic acids research*, 25(17):3389–3402, September 1997.
- [2] Aharon Bar-Hillel, Tomer Hertz, Noam Shental, and Daphna Weinshall. Learning a mahalanobis metric from equivalence constraints. *J. Mach. Learn. Res.*, 6:937–965, 2005.
- [3] ScienceDaily Boston College. Biologists build a better mouse model for cancer research, 2008.
- [4] J. P. Brunet, P. Tamayo, T. R. Golub, and J. P. Mesirov. Metagenes and molecular pattern discovery using matrix factorization. *Proc. Natl. Acad. Sci. U.S.A.*, 101:4164–4169, Mar 2004.
- [5] J. L. Bussiere. Species selection considerations for preclinical toxicology studies for biotherapeutics. *Expert Opin Drug Metab Toxicol*, 4:871–877, Jul 2008.
- [6] E. T. Chan, G. T. Quon, G. Chua, T. Babak, M. Trochesset, R. A. Zirngibl, J. Aubin, M. J. Ratcliffe, A. Wilde, M. Brudno, Q. D. Morris, and T. R. Hughes. Conservation of core gene expression in vertebrate tissues. *J. Biol.*, 8:33, 2009.
- [7] Jesse Davis and Mark Goadrich. The relationship between precision-recall and roc curves. In *ICML '06: Proceedings of the 23rd international conference on Machine learning*, pages 233–240, New York, NY, USA, 2006. ACM.
- [8] Persi Diaconis. *Group representations in probability and statistics*. Institute of Mathematical Statistics Lecture Notes—Monograph Series, 11. Institute of Mathematical Statistics, Hayward, CA, 1988.
- [9] J. Ernst and Z. Bar-Joseph. STEM: a tool for the analysis of short time series gene expression data. *BMC Bioinformatics*, 7:191, 2006.
- [10] W. Fujibuchi, L. Kiseleva, T. Taniguchi, H. Harada, and P. Horton. CellMontage: similar expression profile search server. *Bioinformatics*, 23:3103–3104, Nov 2007.
- [11] T. Hastie, R. Tibshirani, and J. H. Friedman. *The Elements of Statistical Learning*. Springer, corrected edition, July 2009.
- [12] Wassily Hoeffding. A combinatorial central limit theorem. *The Annals of Mathematical Statistics*, 22(4):558–566, 1951.
- [13] L. Hunter, R. C. Taylor, S. M. Leach, and R. Simon. GEST: a gene expression search tool based on a novel Bayesian similarity metric. *Bioinformatics*, 17 Suppl 1:S115–122, 2001.
- [14] Lars J. Jensen, Thomas S. Jensen, Ulrik de Lichtenberg, Søren Brunak, and Peer Bork. Co-evolution of transcriptional and post-translational cell-cycle regulation. *Nature*, September 2006.
- [15] H.S. Le, Z.N. Oltvai, and Z. Bar-Joseph. Cross-species queries of large gene expression databases. *Bioinformatics*, 26(19):2416, 2010.
- [16] D. Lee, O. Redfern, and C. Orengo. Predicting protein function from sequence and structure. *Nature reviews. Molecular cell biology*, 8(12):995–1005, December 2007.
- [17] G. Lelandais, V. Tanty, C. Geneix, C. Etchebest, C. Jacq, and F. Devaux. Genome adaptation to chemical stress: clues from comparative transcriptomics in *Saccharomyces cerevisiae* and *Candida glabrata*. *Genome Biol.*, 9:R164, 2008.

- [18] Y. Lu, X. He, and S. Zhong. Cross-species microarray analysis with the oscar system suggests an *insr*–*pax6*–*nqo1* neuro-protective pathway in aging and alzheimer’s disease. *Nucleic Acids Res*, 35(Web Server issue), July 2007.
- [19] Yong Lu, Peter Huggins, and Ziv Bar-Joseph. Cross species analysis of microarray expression data. *Bioinformatics*, 25(12):1476–1483, June 2009.
- [20] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer Series in Operations Research. Springer, New York, 2006.
- [21] Art B. Owen, Josh Stuart, Kathy Mach, Anne M. Villeneuve, and Stuart Kim. A gene recommender algorithm to identify coexpressed genes in *c. elegans*. *Genome Res.*, 13(8):1828–1837, August 2003.
- [22] N. E. Sharpless and R. A. Depinho. The mighty mouse: genetically engineered mouse models in cancer drug development. *Nat Rev Drug Discov*, 5:741–754, Sep 2006.
- [23] Joshua M. Stuart, Eran Segal, Daphne Koller, and Stuart K. Kim. A gene-coexpression network for global discovery of conserved genetic modules. *Science*, 302(5643):249–255, October 2003.
- [24] Andrew I. Su, Tim Wiltshire, Serge Batalov, Hilmar Lapp, Keith A. Ching, David Block, Jie Zhang, Richard Soden, Mimi Hayakawa, Gabriel Kreiman, Michael P. Cooke, John R. Walker, and John B. Hogenesch. A gene atlas of the mouse and human protein-encoding transcriptomes. *Proceedings of the National Academy of Sciences of the United States of America*, 101(16):6062–6067, April 2004.
- [25] P. Tamayo, D. Scanfeld, B. L. Ebert, M. A. Gillette, C. W. Roberts, and J. P. Mesirov. Metagene projection for cross-platform, cross-species characterization of global transcriptional states. *Proc. Natl. Acad. Sci. U.S.A.*, 104:5959–5964, Apr 2007.
- [26] M. L. Whitfield, G. Sherlock, A. J. Saldanha, J. I. Murray, C. A. Ball, K. E. Alexander, J. C. Matese, C. M. Perou, M. M. Hurt, P. O. Brown, and D. Botstein. Identification of genes periodically expressed in the human cell cycle and their expression in tumors. *Mol. Biol. Cell*, 13:1977–2000, Jun 2002.