

# Meta-analysis of liver transcriptomic data to identify mammalian functional orthologs

Pía Francesca Loren Reyes, Tom Michoel, Anagha Joshi, Guillaume Devailly

*The Roslin Institute, The University of Edinburgh, Easter Bush, Midlothian, EH25 9RG, Scotland, UK*

*Correspondence: Anagha.Joshi@roslin.ed.ac.uk, Guillaume.Devailly@roslin.ed.ac.uk*

---

## Abstract

Identifying which orthologs share functions from sequence alone can be challenging, notably in case of paralogous genes families. We hypothesised that co-expression network will help predict functional orthologs amongst complex homologous gene families. To explore the use of transcriptomic data available in public domain to predict functionally equivalent orthologs, we collected genome wide expression data in mouse and rat liver from over 1500 experiments with varied treatments. We used a hyper-graph clustering method to identify clusters of orthologous genes co-expressed in both mouse and rat. We validated these clusters by analysing expression profiles in each species separately, and demonstrating a high overlap. We then focused on genes in 18 homology groups with one-to-many or many-to-many relationships between two species, to discriminate between functionally equivalent and non-equivalent orthologs.

*Keywords:* gene function, transcriptomics, liver, orthologs, paralogs, co-expression, gene networks

---

## 1. Introduction

Annotation of gene function is a crucial step to understand the DNA sequencing data currently generated at an unprecedented rate. The lack of functional annotation forms a major bottleneck in analyses across diverse fields, including *de novo* genome sequencing [1], Genome Wide Association Studies (GWAS) in model and non-model organisms [2], and metagenomics

[3]. An experimental validation of each gene is impractical to this end as it demands high financial and time cost. It is estimated that only one percent of proteins have experimental functional annotations [4]. Bioinformatic approaches therefore provide an attractive alternative [5]. The most widely used and successful gene annotation strategy has been the annotation transfer between homologous genes. Automated annotation pipelines from sequence alone are widely used, including GOtcha [6] and BlastGO [7]. They allow fast annotation of thousands of genes for newly sequenced genomes [8]. This approach can be used within a species, where gene families (paralogs), might share common functions, or across species, where known function(s) of a gene in one species are used to infer functions of the homologous gene(s) in another species.

Despite being widely used, fast computational annotation comes at a cost of misannotation, which is present at high levels (over 10 percent) and is believed to be increasing [9] due to misannotation transfer. The most common misannotation is over-annotation, where a gene is assigned a specific but incorrect function [10]. This is partly because one of the major challenges in functional annotation transfer across species is that the orthology relationships are not always one-to-one. Specifically, a single gene in one species can be homologous to multiple paralogs in another (one-to-many homologies), after gene duplication or gene loss event(s). After a gene duplication, the two paralogs can have redundant functions, and thus would share similar functional annotations, or one copy might diverge (lose functionality, or gain new functionalities, or change cellular localisation or tissue specificity), and thus paralogs should have different functional annotations despite their homology. Similarly, multigene families (with many-to-many homologies) are highly prone to over-annotation errors.

Protein structure information can act as source for functional distinction within multigene family proteins [4]. Protein-protein interaction networks have also been successfully used to identify functional orthologs [11]; two orthologs interacting with the same proteins in each species are likely to share similar functions. Similar strategy has been applied to biochemical pathway information [12]. We here explore the use of co-expression gene networks in this context, as they offer two main advantages over protein-protein interactions and biochemical pathways. First, they can be inferred from transcriptomic datasets, which are more abundant than protein-protein interaction datasets. Second, they allow functional annotation of the various classes of RNA genes. We have previously shown that multi-species

information improves gene network reconstruction [13].

In order to explore the potential of co-expressed gene networks to identify functional equivalents in complex homologous families, we collected transcriptomic data from mouse and rat liver samples. To minimise technical variation, we collected datasets generated using a single microarray platform in each species, resulting into 920 experiments in mouse and 620 experiments in rat. We firstly identified clusters of co-expressed genes using hierarchical clustering and found biologically relevant clusters. We applied a hyper-graph clustering method, SCHype [14] to simultaneously cluster co-expressed orthologous genes between species. We then focussed on 18 complex (one-to-many or many-to-many) homologous groups, where at least one member in mouse and in rat were present in similar co-regulated gene clusters. This allowed the identification of functional homologs in complex homologous families. Our results show the potential of this method to predict functional orthologs and limit over-zealous annotation transfers.

## 2. Methods

### 2.1. Data collection and normalisation

Microarray data for liver samples in mouse and rat were collected from GEO, where data for mouse was generated using Affymetrix Mouse Genome 430 2.0 Array, and data for rat was generated using Affymetrix Rat Genome 230 2.0 Array as they were the platforms with a large number of experiments available for each species. Experiments came from 62 (mouse) and 28 (rat) independent studies or GEO series. The GEO accession numbers for individual studies are provided in supplementary table 1. Processed data was not directly comparable between studies, as different studies used different normalisation methods, leading to different distribution of values (supplementary figure 1, A and B). As some datasets had a trimmed lower quartile for reduction in noise by limiting the variability of lowly expressed genes, we applied lower quartile trimming on all datasets (supplementary figure 1, C and D). Specifically, we set the expression value of all probes belonging to the lower quartile to the value of the 25 percentile. We then applied quantile normalisation resulting into a uniform distribution of values for each experiment. To facilitate the comparison between mouse and rat data, we used mouse data as a target for quantile normalisation in rat, using preprocessCore functions `normalize.quantiles.determine.target` and `normalize.quantiles.use.target` [15]. Mouse data was selected as

81 the target because it contained more experiments than the rat dataset. Thus,  
82 after our normalisation steps, the distribution of values was identical for each  
83 experiment in both species.

## 84 *2.2. Data clustering*

85 We selected genes with variable expression across experiments by select-  
86 ing probes with a standard deviation greater than one across experiments.  
87 As shown in figure 1, such probes included genes of low as well as high expres-  
88 sion levels, and largely excluded probes showing very low expression in all  
89 experiments. Microarray data being already log-transformed, log fold change  
90 over the average values were obtained by subtracting the mean expression of  
91 each probes.

92 Hierarchical clustering was done on the log fold change matrices using  
93 R functions `dist` and `hclust` with default parameters (euclidean distance,  
94 complete linkage). Dendrogram branches were reordered using the function  
95 `order.optimal` from the `cba` package [16]. Both rows (probes) and columns  
96 (experiments) were clustered using this approach.

97 Gene homology information was retrieved from the Homologen database  
98 [17], and probe orthology information was obtained using the R package an-  
99 notationTools [18]. Due to one-to-many homologs, rat probes and mouse  
100 probes intersections resulted into slightly different numbers for each species.  
101 Average of the two numbers was used to obtain Jaccard indexes. Jaccard  
102 index significance was obtained using the hypergeometric test, and p-values  
103 were corrected for multiple testing using Bonferroni correction. Gene ontol-  
104 ogy enrichment analyses were done using pantherdb [19], using either genes  
105 analysed by the microarray or our list of variable genes as a background.

106 SCHype takes as input a list of conserved interactions which was gen-  
107 erated as follows. First Spearman correlation coefficient between each pair  
108 of probes was obtained independently for both Mouse and Rat expression  
109 data. Pairs of probes with a correlation coefficient greater or equal to 0.5  
110 were selected. Then if orthologs of two connected probes were connected in  
111 the other species, they were kept as an SCHype input. SCHype was run with  
112 a `minClustsize` of 10, otherwise using default parameters. SCHype identi-  
113 fied 132 clusters of homologous genes co-expressed both in mouse and in rat,  
114 which included 825 nodes in mouse and 778 nodes in rat. SCHype allows  
115 probes to be included in multiple clusters. The different number of probes  
116 in mouse and rat is due to the presence of one-to-many and many-to-many

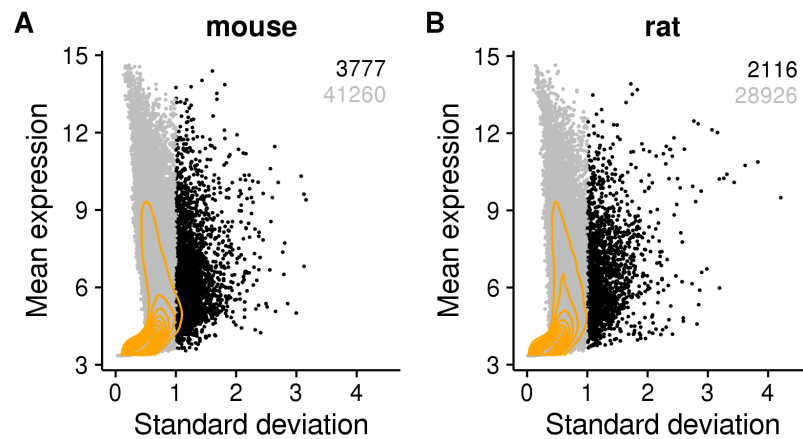


Figure 1: Identification of variable probes in mouse (A) and rat (B) datasets. Each dot represent a single probe. X axis: standard deviation across experiments. Y-axis: mean expression value across experiments (arbitrary units from micro-array processing). In black are the probes with a standard deviation  $\geq 1$ , in grey the probes with a standard deviation  $< 1$ . Orange lines: 2D kernel density.

117 orthologs, as well as the presence of gene measured by multiple probes on  
118 the array.

### 119 2.3. Scripts and data availability

120 R scripts used for this analysis are available in a Github repository [https://github.com/gdevailly/liver\\_mouse\\_rat](https://github.com/gdevailly/liver_mouse_rat). Normalised expression ma-  
121 trices, fold change matrices, as well as probe clusters (hierachical clustering  
122 and SCHype clustering) are available through a Zenodo collection <https://zenodo.org/record/439483>.  
124

## 125 3. Results

### 126 3.1. Identification of variable genes across datasets

127 We downloaded 920 and 620 experiments for gene expression data in  
128 rat and mouse liver from the GEO database. We firstly normalised the  
129 data using lower quartile trimming (supplementary figure 1, C and D) and  
130 quantile normalization (supplementary figure 1, E and F) independently for  
131 each species. We then selected the probes with dynamic expression across  
132 samples (standard deviation  $\geq 1$ ). This resulted into 3777 probes in mouse

Species	Category	Term	Gene	FE	P-value
Mouse	Reactome	Synthesis of (16-20)-hydroxyeicosatetraenoic acids (HETE)	11	4.78	4.29E-02
		Activation of gene expression by SREBF (SREBP)	15	4.34	5.18E-03
		Regulation of cholesterol biosynthesis by SREBP (SREBF)	17	3.94	4.36E-03
		Cytochrome P450 - arranged by substrate type	27	2.72	7.78E-03
	GO slim BP	Phase 1 - Functionalization of compounds	37	2.55	7.58E-04
		fatty acid metabolic process	52	2.26	2.95E-05
		steroid metabolic process	50	2.18	1.31E-04
Rat	Reactome	Synthesis of bile acids and bile salts via 24-hydroxycholesterol	7	8.63	2.95E-02
		Endosomal/Vacuolar pathway	10	7.93	1.15E-03
		Striated Muscle Contraction	11	6.78	1.48E-03
		ER-Phagosome pathway	10	6.53	6.29E-03
		Activation of gene expression by SREBF (SREBP)	10	6.53	6.29E-03
		Antigen Presentation: Folding, assembly and peptide loading of class I MHC	13	6.27	3.84E-04
		Regulation of cholesterol biosynthesis by SREBP (SREBF)	10	5.55	2.51E-02
		Biological oxidations	25	2.95	3.58E-03
	GO slim BP	Metabolism of lipids and lipoproteins	68	2.13	1.15E-05
		response to biotic stimulus	12	4.16	1.12E-02
		fatty acid metabolic process	22	2.52	2.52E-02

Table 1: Variable genes are enriched for categories and pathways relative to liver functions. FE: Fold enrichment between actual over expected number of genes. GO: gene ontology. BP: biological process. Only category with a fold change higher than 2 are shown. All P-values were corrected for multiple testing with the Bonferroni method.

(8.4%) and 2116 probes in rat (6.8%), with a wide range of expression values (figure 1). 735 mouse variable probes out of 3777 had a homologue in rat variable probes, and 624 rat variable probes out of 2116 had a homologue in mouse variable probes. Variable genes were enriched for pathways and functions relative to liver biology (table 1), including metabolism of lipid and protein (rat, adjusted P value  $\leq 10^{-4}$ ), regulation of cholesterol biosynthesis by SREBP (mouse and rat, respectively adjusted P value  $\leq 0.01$  and  $\leq 0.03$ ), synthesis of bile acid and salt via 24-hydroxycholesterol (rat, adjusted P value  $\leq 0.03$ ), and fatty acid metabolic process (mouse and rat, respectively adjusted P value  $\leq 10^{-4}$  and  $\leq 0.03$ ). As the biological processes enriched in variable genes reflected functions associated with liver, we concluded that the expression variability across samples was reflecting biological variability, and not only technical variations, and therefore was of significance for further investigation.

### 147 3.2. Independent hierarchical clustering of mouse and rat data

148 Hierarchical clustering was applied to the mouse and rat expression ma-  
 149 trices independently (figure 2, A and B). We defined 7 major clusters of  
 150 variable probes, while the experiments were grouped in 4 clusters. The two  
 151 major clusters of experiments in mouse showed broadly opposite expression  
 152 patterns (figure 2A). Two major experimental groups were also noted in rat,  
 153 albeit to a lesser extent compared to mouse (figure 2B). Experiments were  
 154 annotated according to their series of origin (figure 2A and B, bottom of  
 155 the heatmap), revealing that most experiments from the same series grouped  
 156 together (including cases and controls). Notably, no series of experiments  
 157 were split in the two main experiment clusters.

158 We characterised the main experiment clusters by looking at the most  
 159 different non-trivial terms in the element-term matrix build from the meta-  
 160 data retrieved from GEO (characteristic field, figure 2, C and D). No clear  
 161 difference between experiment clusters was observed in mouse. Experiment  
 162 cluster 3 in rat seems to be composed mostly of F344 strains of rat and/or of  
 163 rat treated with the microcystinlr toxin. To note, this cluster is dominated  
 164 by experiments from a single GSE (figure 2B). Since experiment clustering  
 165 matched series of origin of the data, this hinders the discrimination of batch  
 166 effects from biological differences.

167 Given that mouse and rat probes formed two major clusters anti-correlated  
 168 with each other despite diverse experimental set ups in each species, we inves-  
 169 tigated whether the mouse and rat probe clusters were composed of probes  
 170 measuring similar genes (figure 2E). We calculated the overlap between genes  
 171 in each cluster in mouse with genes in each cluster in rat. Cluster 2 in mouse  
 172 (golden color, figure 2A) and cluster 2 in rat (golden color, figure 2B) showed  
 173 a very high overlap with the highest Jaccard index across all clusters. Neither  
 174 mouse cluster 2 nor rat cluster 2 were enriched for any gene ontology term  
 175 or reactome pathways, when using the set of variable probes as background.  
 176 Most clusters did not show a very high genes overlap across species, though  
 177 the functional enrichment analysis were suggestive that observed gene vari-  
 178 ations reflected differences in the liver physiology. Specifically, cluster 1 in  
 179 mouse (claret red color, figure 2A) was enriched for generation of precursor  
 180 metabolites and energy (adjusted P value  $\leq 10^{-6}$ ), steroid metabolic process  
 181 (adjusted P value  $\leq 0.001$ ), fatty acid metabolic process (adjusted P value  
 182  $\leq 0.001$ ), and Cytochrome P450 - arranged by substrate type (adjusted P  
 183 value  $\leq 10^{-6}$ ). Cluster 3 in mouse (green color) was enriched for arachi-  
 184 donic acid metabolic process (adjusted P value  $\leq 0.01$ ), icosanoid metabolic



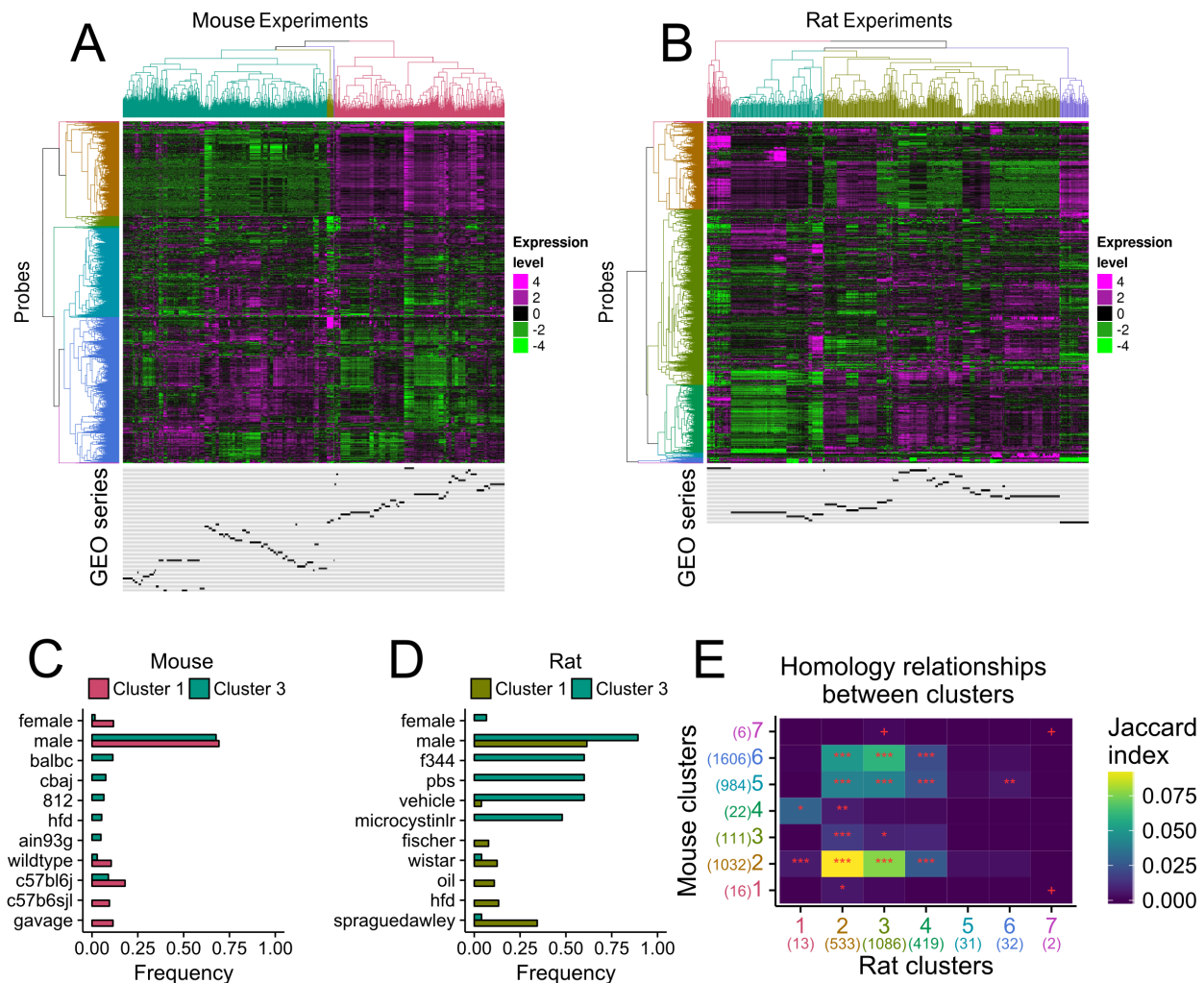


Figure 2: Hierarchical clustering of variable probes in mouse (A) and in rat (B). Four clusters were defined for experiments and seven for probes. Below the heatmaps, localisation of experiments from each series were shown in black, one line per series. FC: fold change. C and D. Metadata term frequencies of the two biggest experiment clusters were compared for mouse (C) and rat (D). Colour-code matches the experiments trees in panels A and B. E. Homology relationship between probe clusters between rat (x-axis) and mouse (y-axis). Cell colour: jaccard index. Cell label: Bonferonni adjusted p-values: \*\*\*  $\leq 0.0001$ , \*\*  $\leq 0.001$ , \*  $\leq 0.01$ , +  $\leq 0.05$ .



process (adjusted P value  $\leq 0.05$ ), fatty acid derivative metabolic process (adjusted P value  $\leq 0.05$ ), and Cytochrome P450 - arranged by substrate type (adjusted P value  $\leq 0.05$ ). Cluster 6 in rat (blue) was enriched for proteolysis (FE 10, adjusted P value  $\leq 0.01$ ). More term relative to the liver metabolism were enriched when the same analysis was used using all genes as a background (supplementary table 2).

### 3.3. Co-clustering of Mouse and Rat expression data

To identify clusters of homologous probes between mouse and rat, we used the hyper-graph clustering tool SCHype [14]. SCHype uses a recursive spectral clustering algorithm to identify sets of nodes in each species with a greater than expected number of conserved interactions between them (figure 3A). Input data for SCHype was built using three graphs: a mouse probe graph built from pairs of probes with a Spearman correlation coefficient  $\geq 0.5$  (supplementary figure 2A), a rat probe graph with pairs of probes with a Spearman correlation coefficient  $\geq 0.5$  (supplementary figure 2B), and a probe to probe homology graph between rat and mouse built using the Homologene database [17] and the annotationTools package [18]. SCHype identified 132 clusters of homologous genes co-expressed both in mouse and in rat, which included 825 nodes in mouse and 778 nodes in rat (figure 3B). SCHype allows probes to be included in multiple clusters resulting into 474 unique probes in mouse and 425 unique probes in rat. It identified four clusters with over 30 homologous genes in each species, eighteen clusters with over 10 probes in each species, thirty-five clusters with only 2 co-expressed probes in each species (figure 3B). We further focussed on the first four (c1-c4) SCHype clusters (figure 3C). We firstly compared SCHype clusters with results obtained by clustering data from each species independently. SCHype cluster c3 highly overlapped with the previous cluster 2 in mouse (golden color, figure 2A) and the cluster 2 in rat (golden color, figure 2B). These two clusters were shown to share a high number of homologous probes (figure 2E). Gene ontology analysis of the four biggest SCHype clusters, both over the set of variable probes or over the full set of probes, did not lead to any significant results, most likely due to small number of genes in each cluster. Importantly, the series experiments in each series no longer clustered together after restricting the data to each of the four biggest SCHype clusters (figure 3C). Individual experiments from each series nevertheless belonged to the same large experiment cluster (figure 3C) highlighting the need for building an expression compendium to obtain these results.

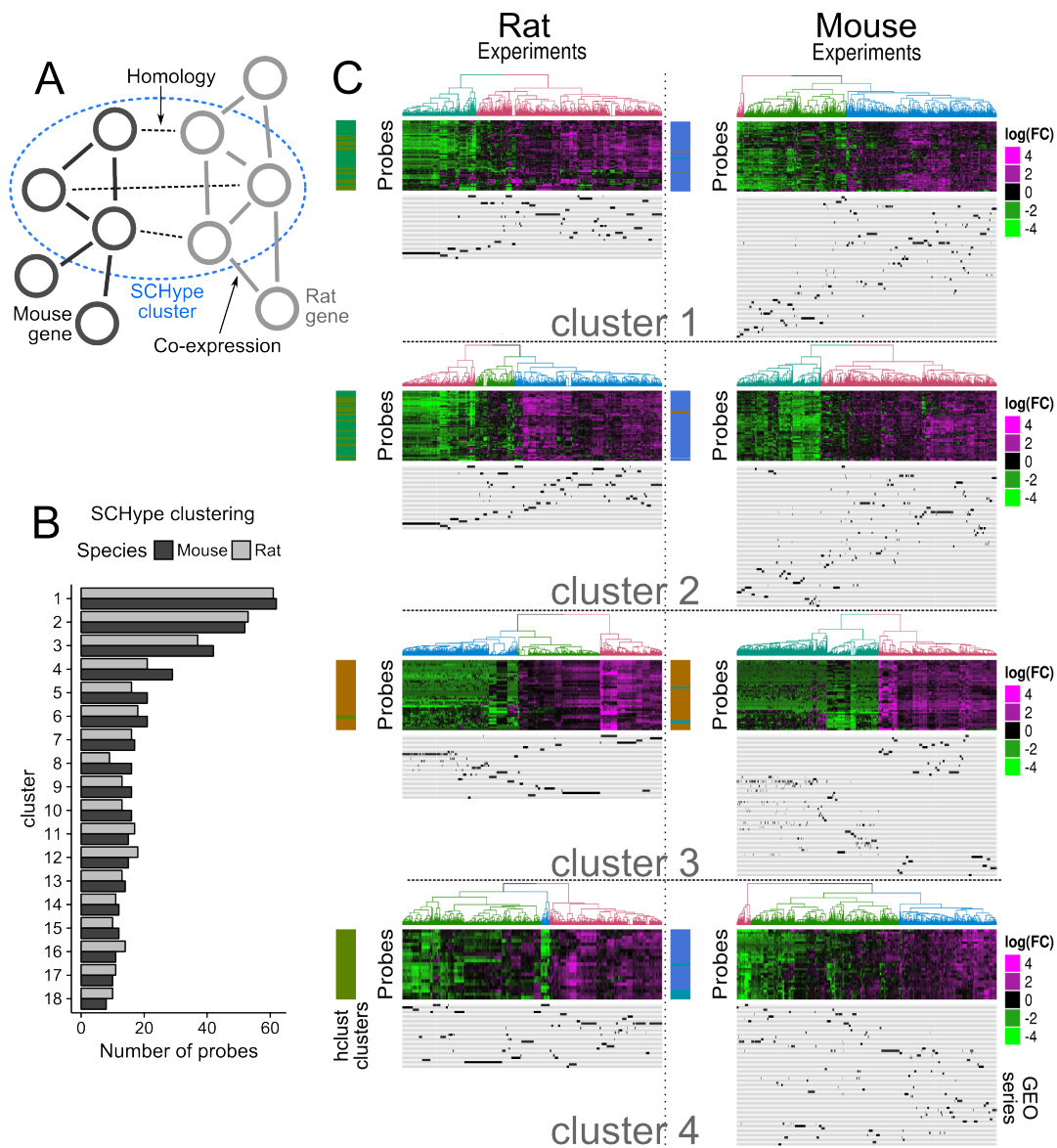


Figure 3: Co-clustering of rat (middle) and mouse (right) data using SCHype. A. SCHype is a clustering tool for hypergraphs, built here from two co-expression graphs and an homology graph. B. Number of mouse (dark grey) and rat (light grey) probes for the SCHype clusters with more than 10 probes for each species. C. The biggest four SCHype clusters are shown. Genes in mouse and rat in each cluster are homologous to each other. The results of hierarchical clustering for each species is shown as a color bar on the left. Colour-code matches the experiments trees in figure 1. Under the heatmap, clustering localisation of experiments from each series are shown in black, one line per series.

Homology group	Species	Gene name	SCHype cluster	
137229			cluster 69	
	mouse	<i>Anp32a</i>	✓	
	rat	<i>Anp32a</i>	✓	
	rat	<i>LOC100909983</i>		
68982			cluster 7	cluster 30
	mouse	<i>Ccnb1</i>	✓	✓
	mouse	<i>Gm5593</i>		
	rat	<i>Ccnb1</i>	✓	✓
10699			cluster 2	cluster 118
	mouse	<i>Cd248</i>	✓	✓
	rat	<i>Cd248</i>	✓	✓
	rat	<i>LOC100911932</i>		
	rat	<i>LOC100911882</i>		
3938			cluster 1	
	mouse	<i>Ppp1r3c</i>	✓	
	rat	<i>Ppp1r3c</i>	✓	
	rat	<i>LOC100910671</i>		
14108			cluster 2	
	mouse	<i>Rasl10b</i>	✓	
	rat	<i>Rasl10b</i>	✓	
	rat	<i>LOC100912246</i>		

Table 2: SCHype clustering of homologous groups: predicted functionality is following gene names. Homology group is from the Homologene database. Tick mark indicates the inclusion of the gene in the corresponding SCHype cluster.

### 3.4. Co-clustering across species identifies functionally equivalent orthologs

SCHype clustering successfully identified clusters of homologous genes co-expressed in both mouse and rat datasets. This information adds an independent evidence in support of a functional annotation transfer for pairs of orthologous genes across species found in the same SCHype cluster(s), as functionally equivalent orthologs would be co-expressed with the same set of genes in both species, and therefore would be included in the same SCHype cluster(s). We investigated if SCHype clusters could help identify functionally equivalent orthologs amongst complex homology groups. Eighteen homology groups of three members or more had at least one member of each species in the same SCHype cluster(s). For example, for homology group 137299 (table 2), *Anp32a* in mouse and *Anp32a* in rat were in the same SCHype cluster 69, while *LOC100909983*, another homologue of rat *Anp32a*, was not. This suggests that indeed *Anp32a* in rat is the functional equivalent of *Anp32a* in mouse, but *LOC100909983* is not. In this case, our method found back a functional equivalent already known [20]. Similar observations were made for homology groups 68982 (*Ccnb1*), 10699 (*Cdc248*), 3938 (*Ppp1r3c*), and 14108 (*Rasl10b*) (table 2). In five cases, all members of

Homology group	Species	Gene name	SCHype cluster		
128630			cluster 9	cluster 12	cluster 45
	mouse	<i>Ceacam1</i>	✓		
	mouse	<i>Ceacam2</i>	✓	✓	✓
	rat	<i>Ceacam1</i>	✓	✓	✓
11456			cluster 5		
	mouse	<i>Elovl6</i>	✓		
	rat	<i>Elovl6</i>	✓		
	rat	<i>LOC102549542</i>	✓		
20277			cluster 35		
	mouse	<i>Rrm2</i>	✓		
	rat	<i>Rrm2</i>	✓		
	rat	<i>LOC100359539</i>	✓		
55991			cluster 1	cluster 119	
	mouse	<i>Tmed2</i>	✓	✓	
	mouse	<i>Gm21540</i>	✓	✓	
	rat	<i>Tmed2</i>	✓	✓	
11890			cluster 10	cluster 43	cluster 81
	mouse	<i>Tnks2</i>	✓	✓	✓
	rat	<i>LOC100910717</i>	✓	✓	✓
	rat	<i>Tnks2</i>	✓	✓	✓

Table 3: SCHype clustering of homologous groups: all member of the homology groups share predicted functionalities. Homology group is from the Homologene database. Tick mark indicates the inclusion of the gene in the corresponding SCHype cluster.

the homology groups were included in the same SCHype clusters (table 3), suggesting that all orthologs are likely to share the same function(s). Finally, eight homology groups showed more complex situations, where neither only one nor all the homologs were present in the same groups (table 4 and supplementary table 3). For example, in homology group 117945, *Cyp2c7* in rat had three homologous genes in mouse but only *Cyp2c38* in mouse belonged to the same SCHype cluster (table 4) predicting that mouse *Cyp2c38* (and not mouse *Cyp2c29* or mouse *Cyp2c39*) is a functional ortholog of rat *Cyp2c7*.

## 4. Discussion

Here we have shown that transcriptomic data can be used to help predict functionally orthologous genes, using co-expression networks built from mouse and rat liver samples. We identified 18 complex homologous groups (i.e. with paralogs in at least one of the species), including 54 genes in mouse and 46 genes in rat, with at least one gene in mouse and one gene in rat in the same SCHype cluster(s). Lowering the correlation threshold and the standard deviation threshold will likely increase the number of homologous groups, potentially with a higher false positive rate. In this study, we focussed on a single tissue as a proof of concept. Addition of data from various

Homology group	Species	Gene name	SCHype cluster
117948			cluster 102
	mouse	<i>Cyp2c38</i>	✓
	mouse	<i>Cyp2c29</i>	
	mouse	<i>Cyp2c39</i>	
	rat	<i>Cyp2c7</i>	✓
104115			cluster 33
	mouse	<i>Hsd3b5</i>	✓
	mouse	<i>Gm10681</i>	
	mouse	<i>Hsd3b4</i>	
	mouse	<i>Gm4450</i>	
	rat	<i>Hsd3b5</i>	✓
	rat	<i>LOC100911116</i>	✓
137425			cluster 2
	mouse	<i>Lce3c</i>	✓
	rat	<i>LOC100361951</i>	✓
	rat	<i>LOC100911982</i>	✓
	rat	<i>Lce3d</i>	
129514			cluster 17
	mouse	<i>Rdh9</i>	✓
	mouse	<i>Rdh1</i>	
	mouse	<i>Rdh16</i>	
	mouse	<i>Rdh19</i>	
	mouse	<i>BC089597</i>	
	rat	<i>Rdh16</i>	✓
	rat	<i>LOC100365958</i>	✓

Table 4: SCHype clustering of homologous groups: non-trivial functional relations. Homology group is from the Homologene database. Tick mark indicates the inclusion of the gene in the corresponding SCHype cluster. Four additional, more complex, homology groups are shown in supplementary table 3.

other tissues might lead to co-expressed gene networks of many more genes, but they may lack the fine resolution that is needed to improve functional annotation inference due to lack of a tissue specific context. We used microarray data in this study as it is by far the most abundant dataset. However, consortia like GTEx [21] have generated large amount of RNA sequencing data, and we envisage application of the method described here to RNA sequencing data in the future. The greater sensitivity of RNA sequencing over microarray [22] might allow the identification of more co-expressed genes.

Despite rigorous data normalisation, experiments from the same series tended to cluster together, cases and controls included. While this could be a sign of technical biases, gene ontology analysis of the variable genes demonstrated that they are related to liver functions. Thus it appears that the gene expression variability we observed is, at least partially, reflecting biological variations impacting the liver physiology. Importantly, individual experiments from series did not cluster together in SCHype clusters. We applied various approaches but could not identify the biological origin(s) of the observed variations. This is in part due to the lack of standardised experiment metadata fields in GEO (not all datasets even had a strain or a sex information, for example), and the lack of controlled vocabulary used to describe experiments. It is a possibility that better annotation of the metadata would have allowed the identification of critical confounding factors.

SCHype clustering was able to find some known as well as some yet to be experimentally validated orthologous functional relationships. For example, only mouse and rat *Ccnb1* where in the same SCHype cluster, and not *Gm5593*. While mouse *Ccnb1* and rat *Ccnb1* are annotated as protein coding genes, *Gm5593* in mouse is annotated as a processed pseudogene [20].

Finally, we note that co-expression of orthologous genes is not a validation of shared functionality, but it can be used as another source of evidence. While protein-protein interaction networks could be used for the same aim, transcriptomic data are more easily generated and therefore widely available for many species. Thus the method described here shows a promise to enhance functional gene annotation transfer across species. It can provide an experimental support for one-to-one ortholog annotation transfer, and can help identify functionally similar and non similar orthologs in one-to-many and many-to-many orthology groups.



## 5. Funding

AJ is a Chancellor's fellow and AJ and TM labs are supported by institute strategic funding from Biotechnology and Biological Sciences Research Council (BBSRC, BB/J004235/1). GD is funded by the People Programme (Marie Curie Actions FP7/2007-2013) under REA grant agreement No PCOFUND-GA-2012-600181.

## 6. References

- [1] R. Salavati, H. S. Najafabadi, Sequence-based functional annotation: what if most of the genes are unique to a genome?, *Trends in parasitology* 26 (2010) 225–9.
- [2] S. L. Edwards, J. Beesley, J. D. French, A. M. Dunning, Beyond GWASs: Illuminating the Dark Road from Association to Function, *The American Journal of Human Genetics* 93 (2013) 779–797.
- [3] C. Rinke, P. Schwientek, A. Sczyrba, N. N. Ivanova, I. J. Anderson, J.-F. Cheng, A. Darling, S. Malfatti, B. K. Swan, E. A. Gies, J. A. Dodsworth, B. P. Hedlund, G. Tsiamis, S. M. Sievert, W.-T. Liu, J. A. Eisen, S. J. Hallam, N. C. Kyrpides, R. Stepanauskas, E. M. Rubin, P. Hugenholtz, T. Woyke, Insights into the phylogeny and coding potential of microbial dark matter, *Nature* 499 (2013) 431–437.
- [4] S. Das, C. A. Orengo, Protein function annotation using protein domain family resources., *Methods (San Diego, Calif.)* 93 (2016) 24–34.
- [5] D. Davidson, R. Baldock, Bioinformatics beyond sequence: mapping gene function in the embryo., *Nature reviews. Genetics* 2 (2001) 409–17.
- [6] D. M. A. Martin, M. Berriman, G. J. Barton, GOtcha: a new method for prediction of protein function assessed by the annotation of seven genomes., *BMC bioinformatics* 5 (2004) 178.
- [7] S. Götz, J. M. García-Gómez, J. Terol, T. D. Williams, S. H. Nagaraj, M. J. Nueda, M. Robles, M. Talón, J. Dopazo, A. Conesa, High-throughput functional annotation and data mining with the Blast2GO suite., *Nucleic acids research* 36 (2008) 3420–35.

- 324 [8] D. Cozzetto, D. T. Jones, Computational Methods for Annotation  
325 Transfers from Sequence., *Methods in molecular biology* (Clifton, N.J.)  
326 1446 (2017) 55–67.
- 327 [9] A. M. Schnoes, S. D. Brown, I. Dodevski, P. C. Babbitt, Annotation er-  
328 ror in public databases: misannotation of molecular function in enzyme  
329 superfamilies., *PLoS computational biology* 5 (2009) e1000605.
- 330 [10] R. Zallot, K. Harrison, B. Kolaczowski, V. de Crécy-Lagard, Functional  
331 Annotations of Paralogs: A Blessing and a Curse, *Life* 6 (2016) 39.
- 332 [11] M. Kolár, M. Lässig, J. Berg, From protein interactions to functional  
333 annotation: graph alignment in Herpes., *BMC systems biology* 2 (2008)  
334 90.
- 335 [12] F. Mao, Z. Su, V. Olman, P. Dam, Z. Liu, Y. Xu, Mapping of or-  
336 thologous genes in the context of biological pathways: An application of  
337 integer programming., *Proceedings of the National Academy of Sciences*  
338 of the United States of America 103 (2006) 129–34.
- 339 [13] A. Joshi, Y. Beck, T. Michoel, Multi-Species Network Inference Im-  
340 proves Gene Regulatory Network Reconstruction for Early Embryonic  
341 Development in *Drosophila*, *Journal of Computational Biology* 22 (2015)  
342 253–265.
- 343 [14] T. Michoel, B. Nachtergaele, Alignment and integration of complex  
344 networks by hypergraph-based spectral clustering, *Physical Review E*  
345 86 (2012) 056111.
- 346 [15] B. Bolstad, preprocessCore: A collection of pre-processing functions.,  
347 2016.
- 348 [16] C. Buchta, M. Hahsler, cba: Clustering for Business Analytics, 2017. R  
349 package version 0.2-18.
- 350 [17] Database resources of the National Center for Biotechnology Informa-  
351 tion, *Nucleic Acids Research* 44 (2016) D7–D19.
- 352 [18] A. Kuhn, R. Luthi-Carter, M. Delorenzi, Cross-species and cross-  
353 platform gene expression studies with the Bioconductor-compliant R  
354 package 'annotationTools', *BMC Bioinformatics* 9 (2008) 26.

- 355 [19] H. Mi, X. Huang, A. Muruganujan, H. Tang, C. Mills, D. Kang, P. D.  
356 Thomas, PANTHER version 11: expanded annotation data from Gene  
357 Ontology and Reactome pathways, and data analysis tool enhancements,  
358 Nucleic Acids Research 45 (2017) D183–D189.
- 359 [20] A. Yates, W. Akanni, M. R. Amode, D. Barrell, K. Billis, D. Carvalho-  
360 Silva, C. Cummins, P. Clapham, S. Fitzgerald, L. Gil, C. G. Girón,  
361 L. Gordon, T. Hourlier, S. E. Hunt, S. H. Janacek, N. Johnson, T. Juet-  
362 temann, S. Keenan, I. Lavidas, F. J. Martin, T. Maurel, W. McLaren,  
363 D. N. Murphy, R. Nag, M. Nuhn, A. Parker, M. Patricio, M. Pignatelli,  
364 M. Rahtz, H. S. Riat, D. Sheppard, K. Taylor, A. Thormann, A. Vullo,  
365 S. P. Wilder, A. Zadissa, E. Birney, J. Harrow, M. Muffato, E. Perry,  
366 M. Ruffier, G. Spudich, S. J. Trevanion, F. Cunningham, B. L. Aken,  
367 D. R. Zerbino, P. Flicek, Ensembl 2016., Nucleic acids research 44  
368 (2016) D710–6.
- 369 [21] M. Melé, P. G. Ferreira, F. Reverter, D. S. DeLuca, J. Monlong, M. Sam-  
370 meth, T. R. Young, J. M. Goldmann, D. D. Pervouchine, T. J. Sullivan,  
371 R. Johnson, A. V. Segrè, S. Djebali, A. Niarchou, T. G. Consortium,  
372 F. A. Wright, T. Lappalainen, M. Calvo, G. Getz, E. T. Dermitzakis,  
373 K. G. Ardlie, R. Guigó, The human transcriptome across tissues and  
374 individuals, Science 348 (2015) 660–665.
- 375 [22] Z. Wang, M. Gerstein, M. Snyder, Rna-seq: a revolutionary tool for  
376 transcriptomics, Nature reviews. Genetics 10 (2009) 5763.