

# **Revising the ortholog conjecture in cross-species comparison of single-cell transcriptomics**

Yuyao Song <sup>\*1</sup>, Detlev Arendt <sup>2</sup>, Irene Papatheodorou <sup>\*1,3,4</sup>, Alvis Brazma <sup>1,5</sup>

1 European Molecular Biology Laboratory-European Bioinformatics Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, CB10 1SA, UK

2 European Molecular Biology Laboratory (EMBL), Developmental Biology Unit, Heidelberg 69117, Germany

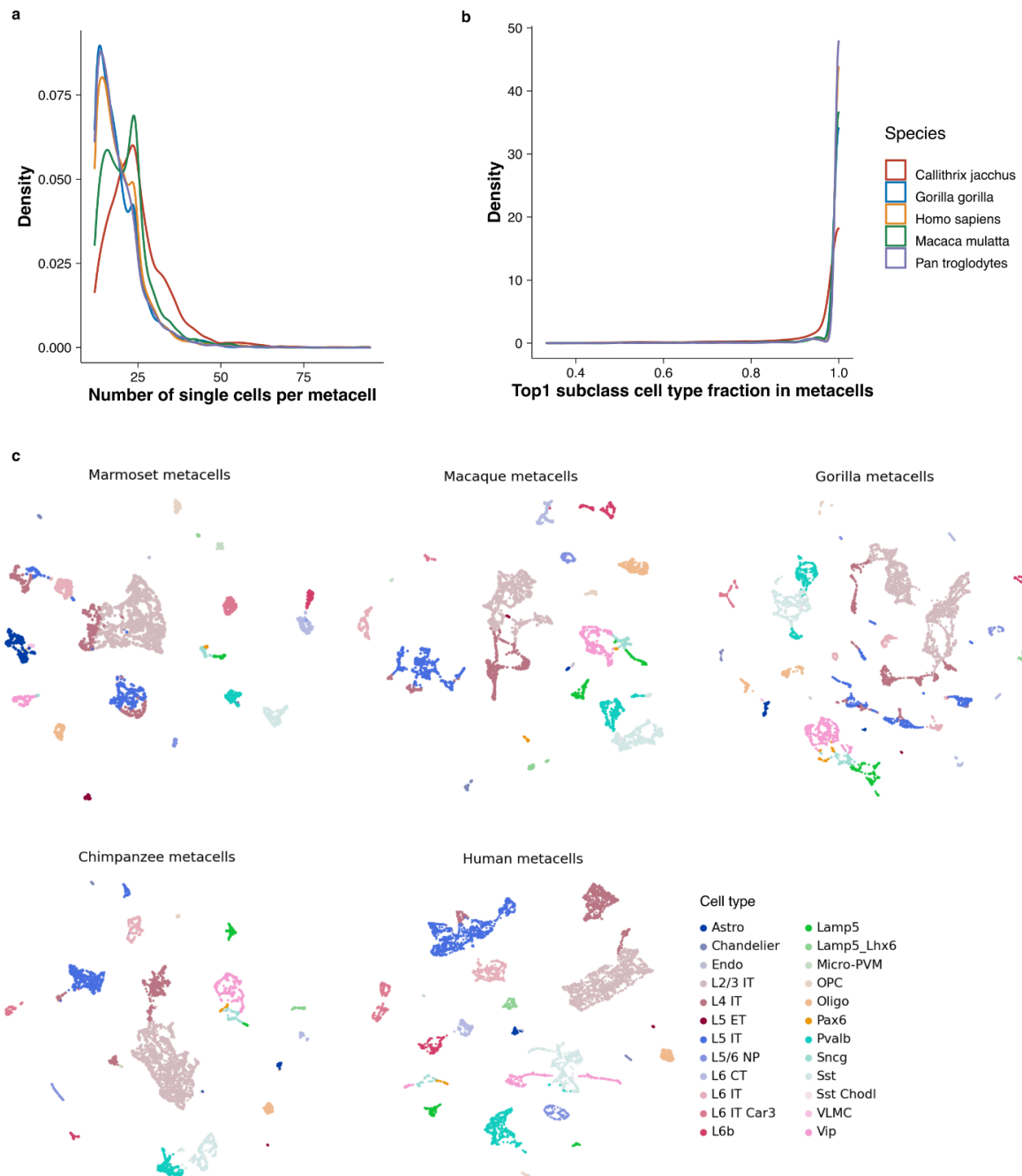
3 Earlham Institute, Norwich Research Park, Norwich, NR4 7UZ, UK

4 Medical School, University of East Anglia, Norwich Research Park, Norwich, NR4 7UA, UK

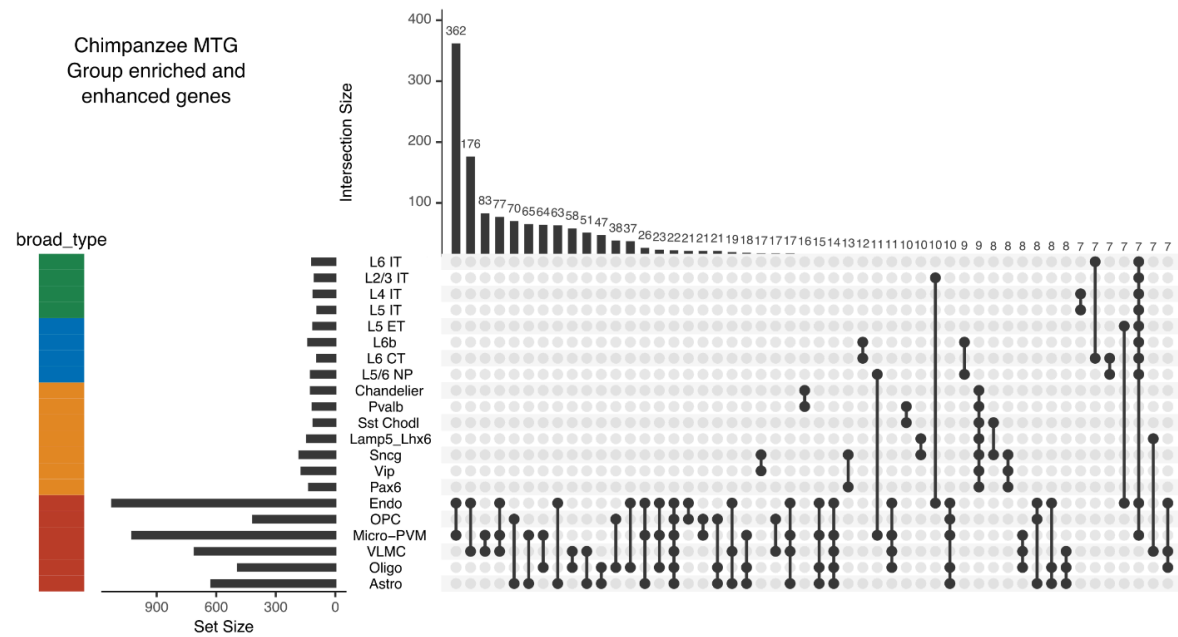
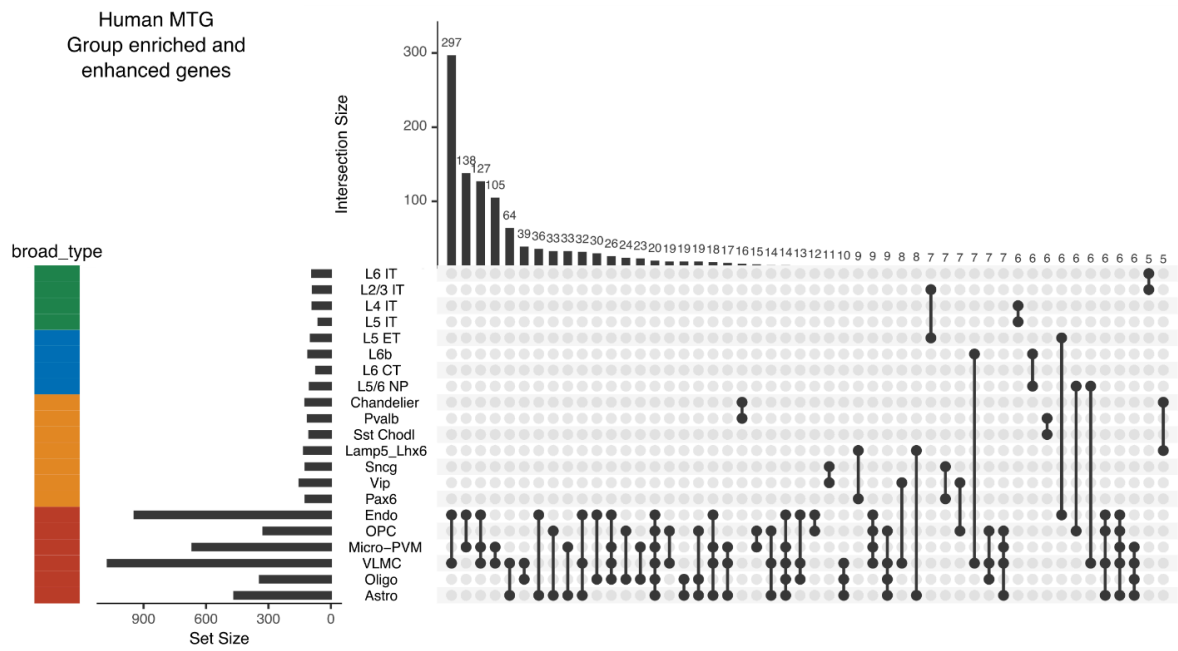
5 Biomedical Research and Study Centre, Ratsupites 1, Riga, LV-1067, Latvia

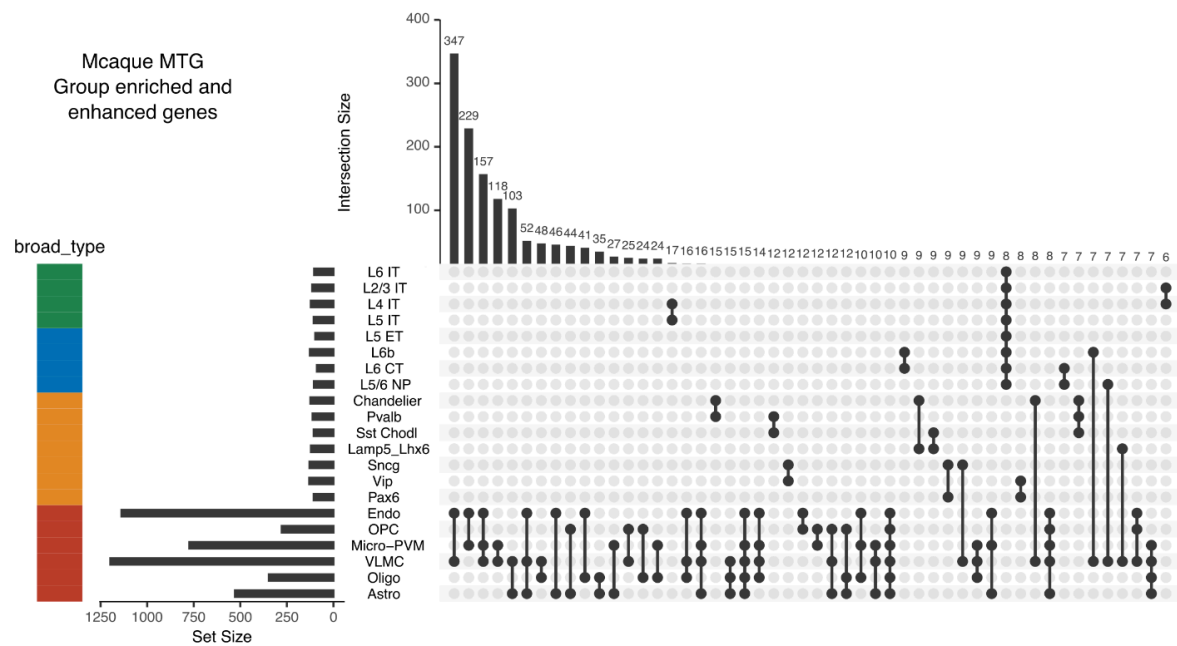
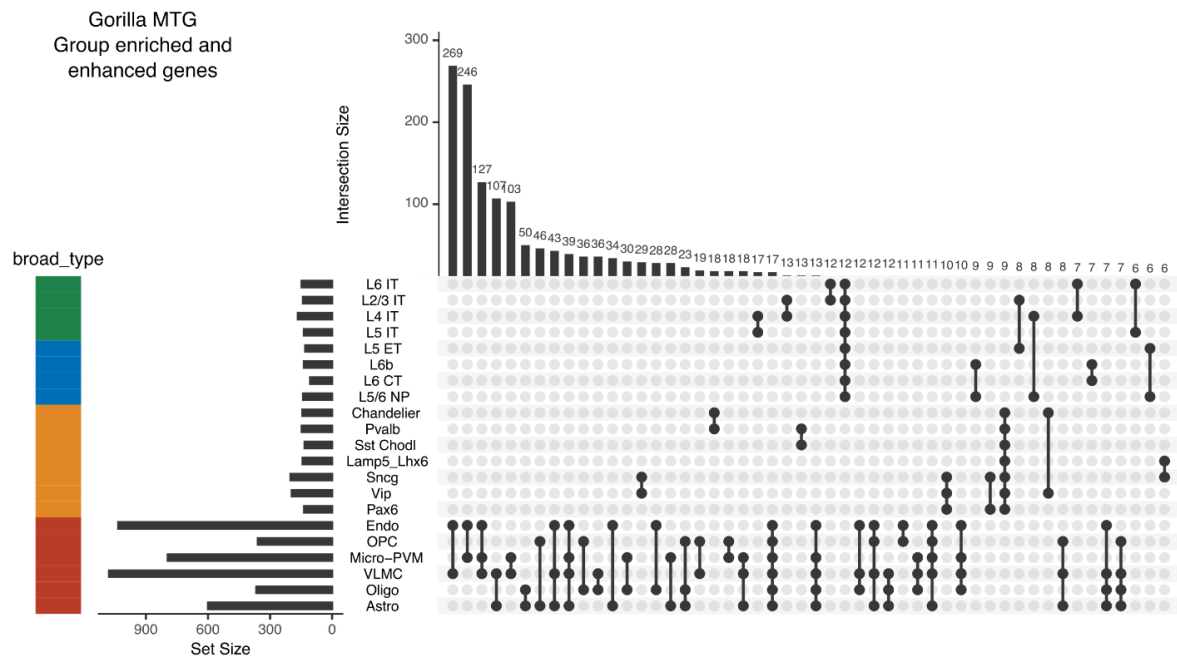
\* Corresponding authors

**Supplementary Figures 1-16**

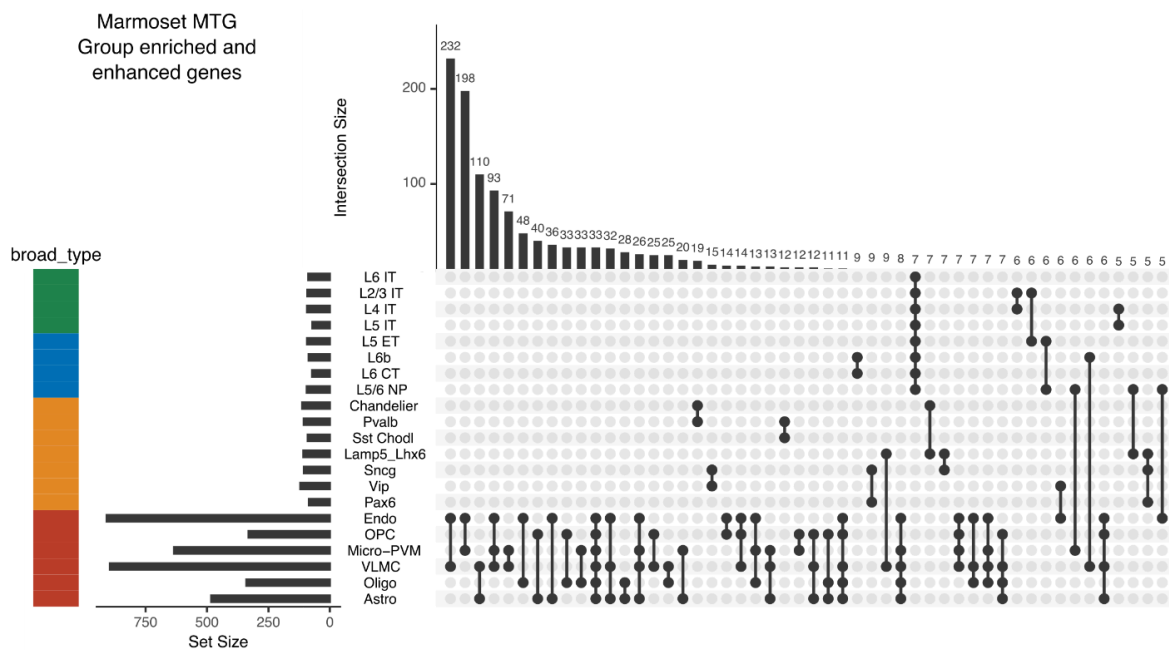


**Supplementary Figure 1 Metacells in the primate MTG data.** (a) Number of single cells per metacell in each species. (2) Fraction of the top 1 cell type in each metacell in each species. (c) UMAP visualisation of metacells in each species with cell type annotation ("Subclass" in original study). MTG, middle temporal gyrus.

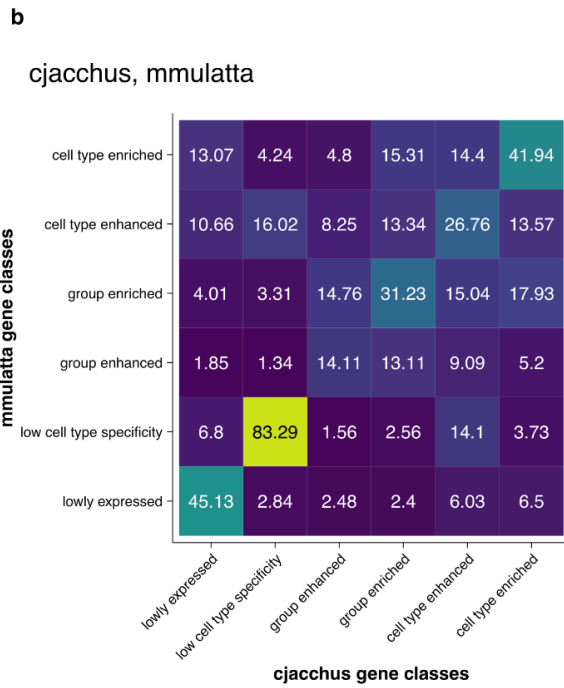
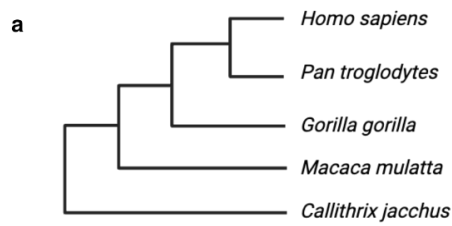




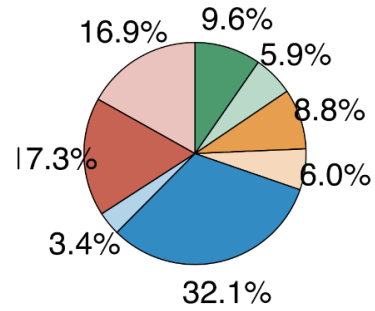




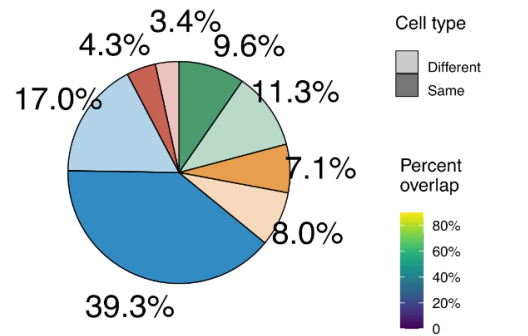
**Supplementary Figure 2 Group specific genes shared among cell types.** Upset plots showing which cell types share the most group specific (enriched and enhanced) genes in each species. Intersections are ordered by intersection size, while set size is ordered by broad type.



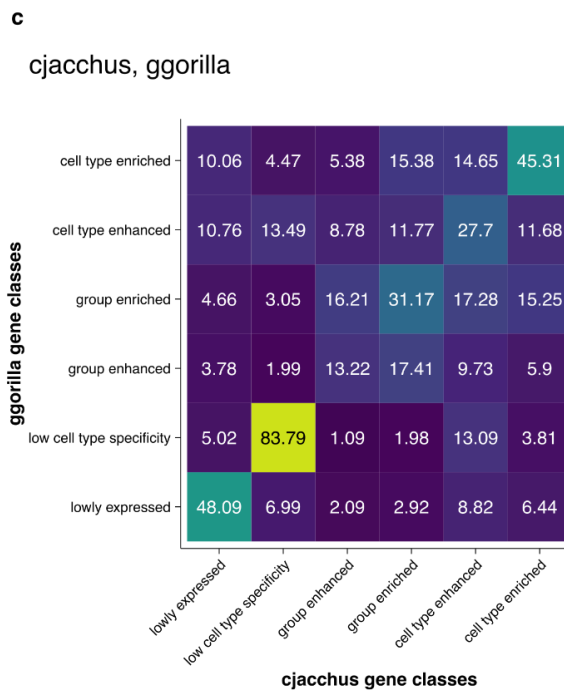
cell type enriched / enhanced genes



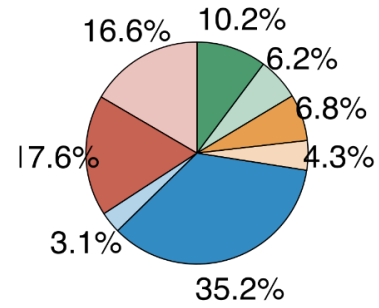
group enriched / enhanced genes



**a**

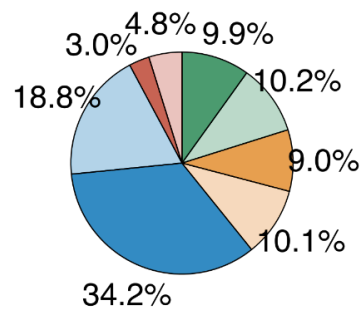


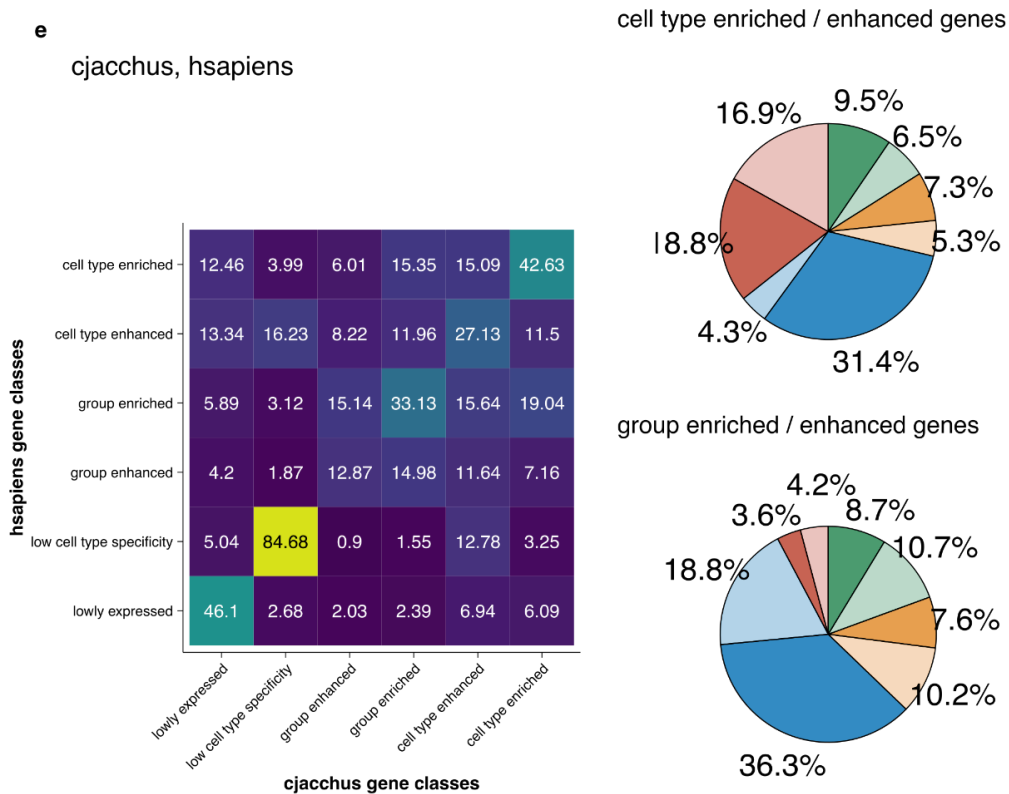
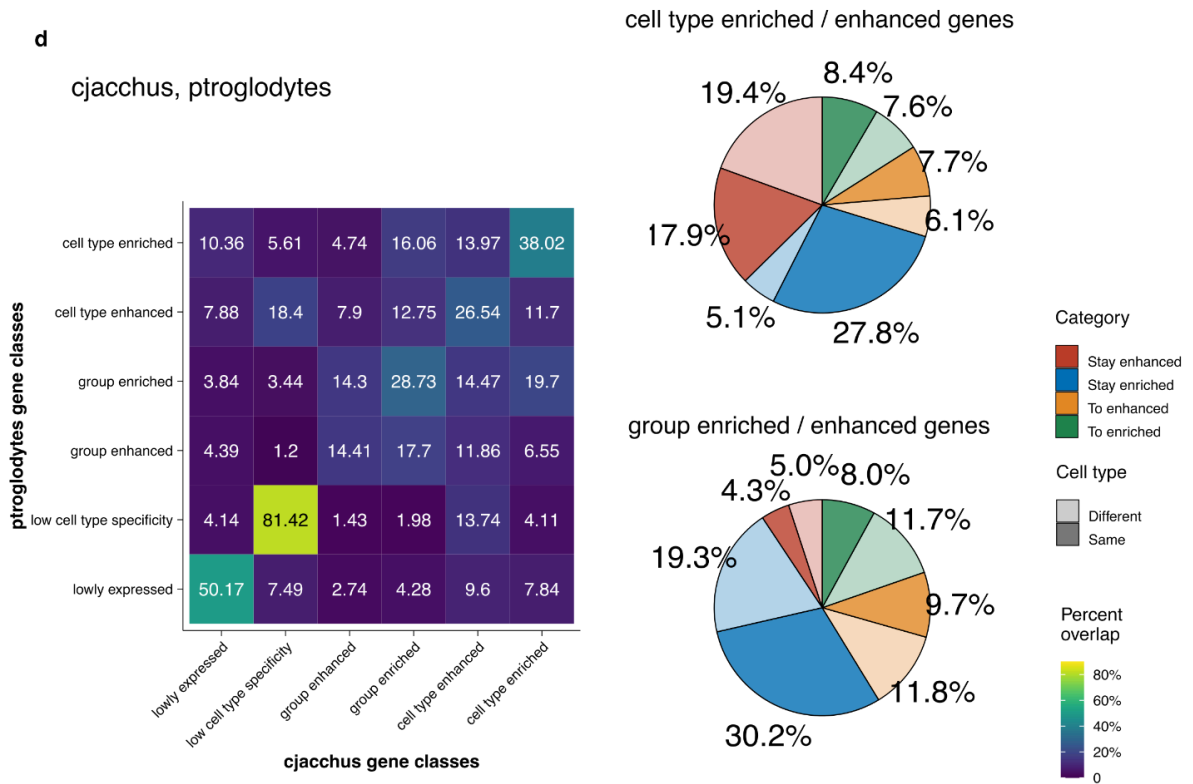
cell type enriched / enhanced genes



**a**

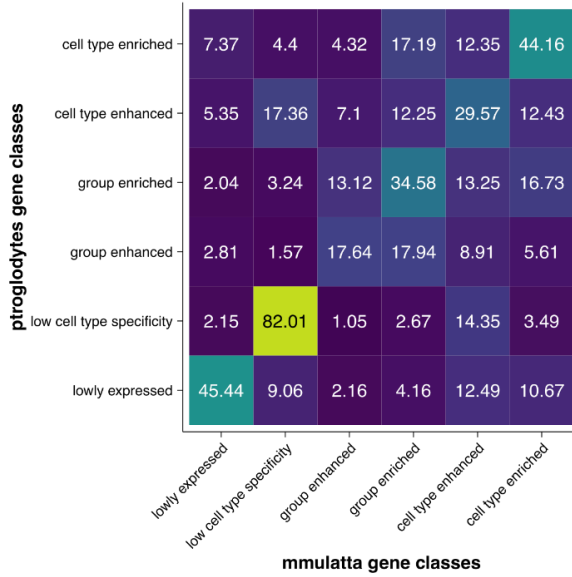
group enriched / enhanced genes



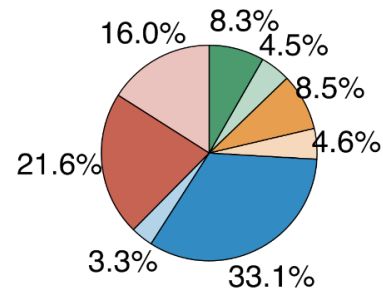


f

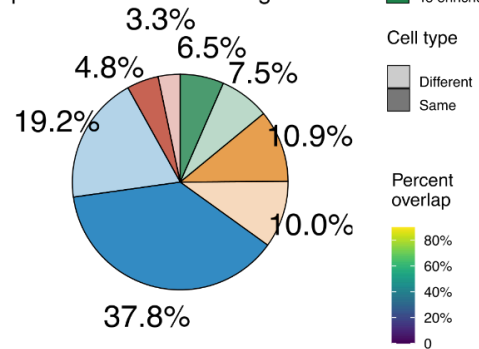
mmulatta\_ptroglodytes



cell type enriched / enhanced genes

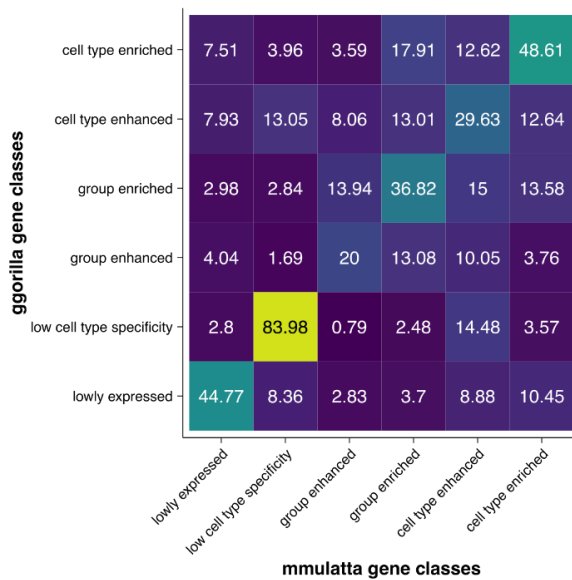


group enriched / enhanced genes

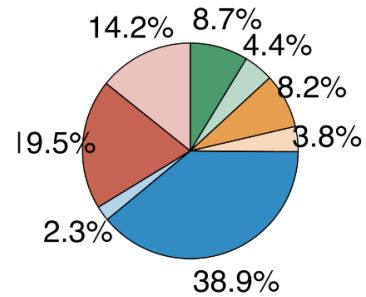


g

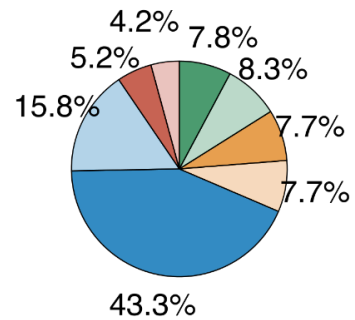
mmulatta\_ggorilla



cell type enriched / enhanced genes

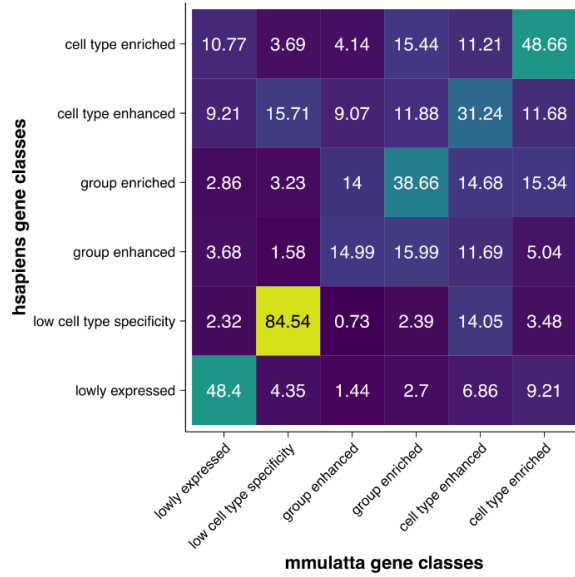


group enriched / enhanced genes

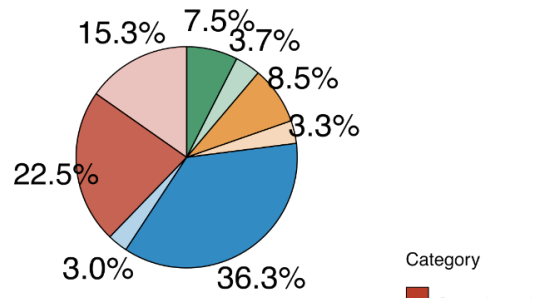


h

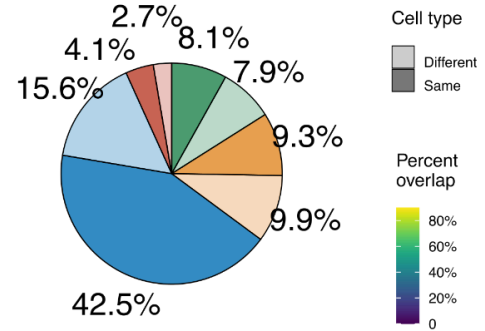
mmulatta\_hsapiens



cell type enriched / enhanced genes

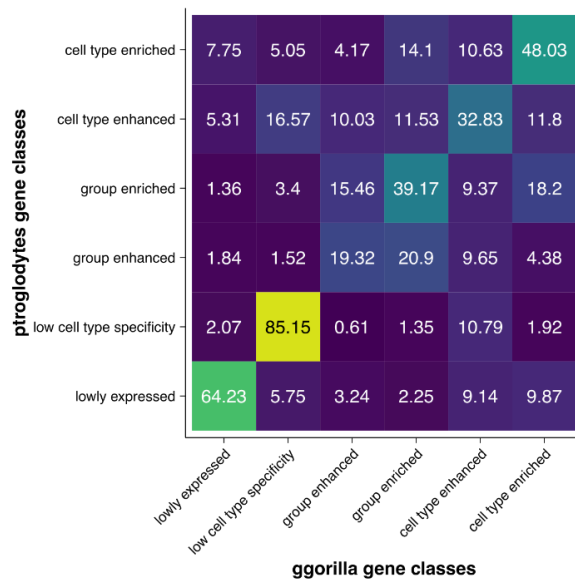


group enriched / enhanced genes

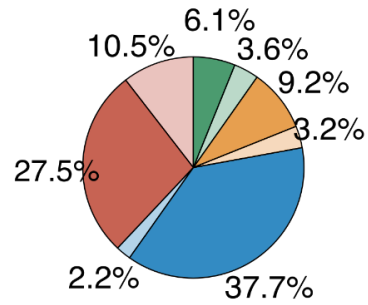


i

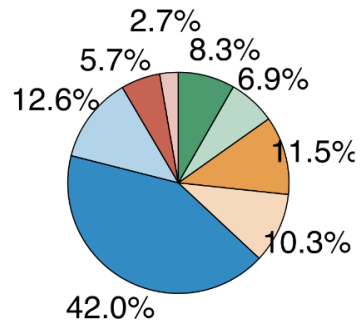
ggorilla\_ptroglodytes

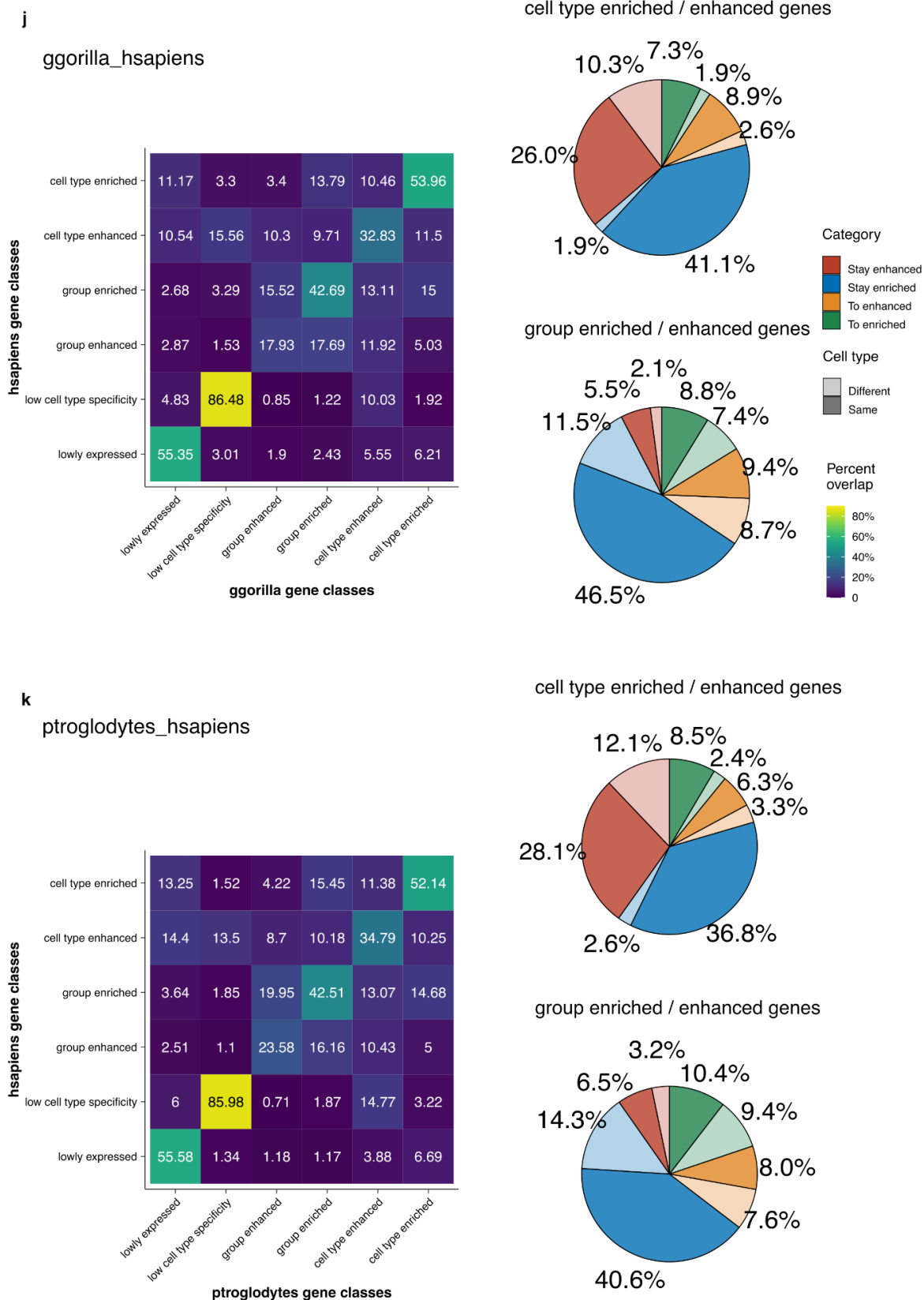


cell type enriched / enhanced genes



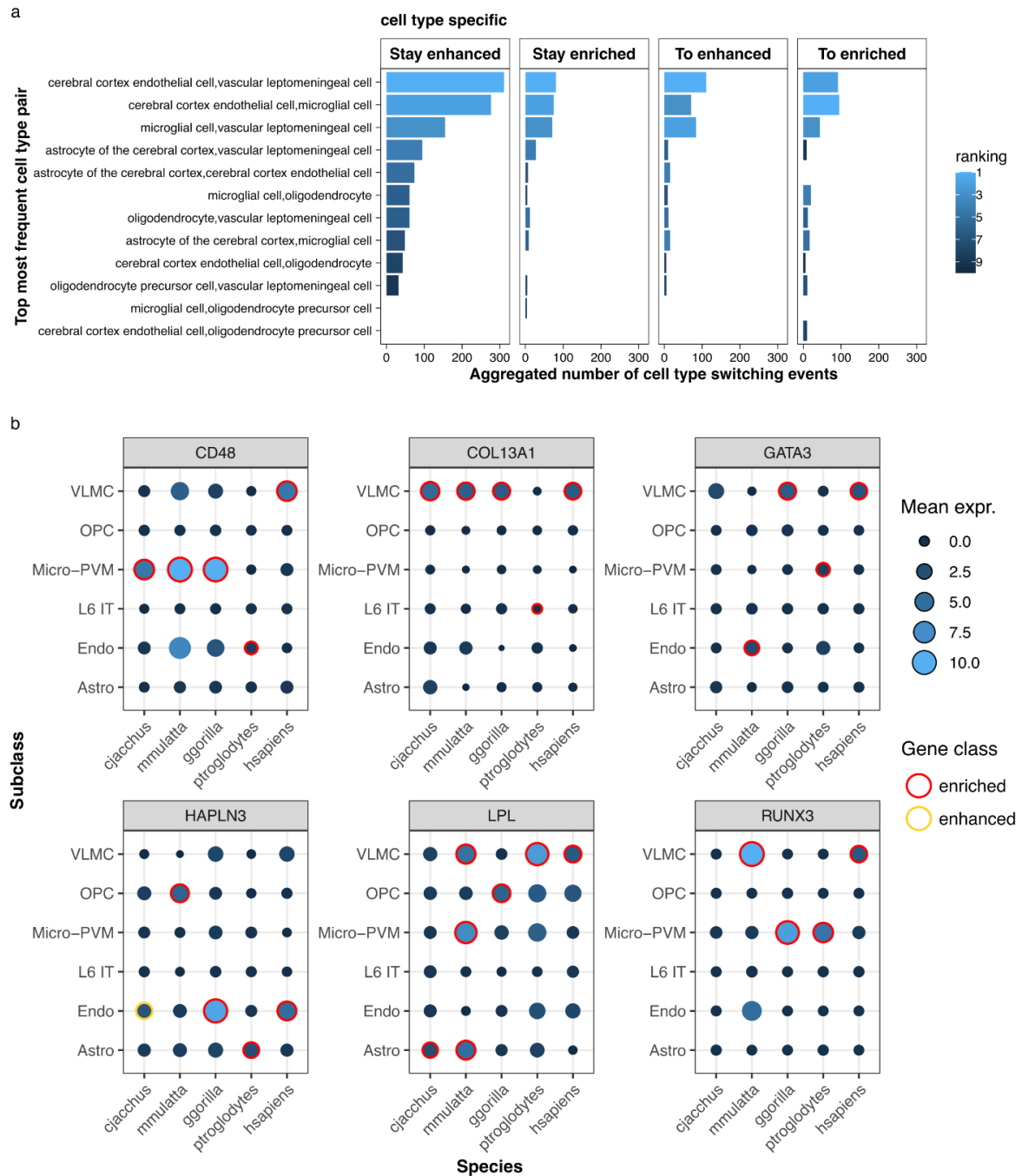
group enriched / enhanced genes





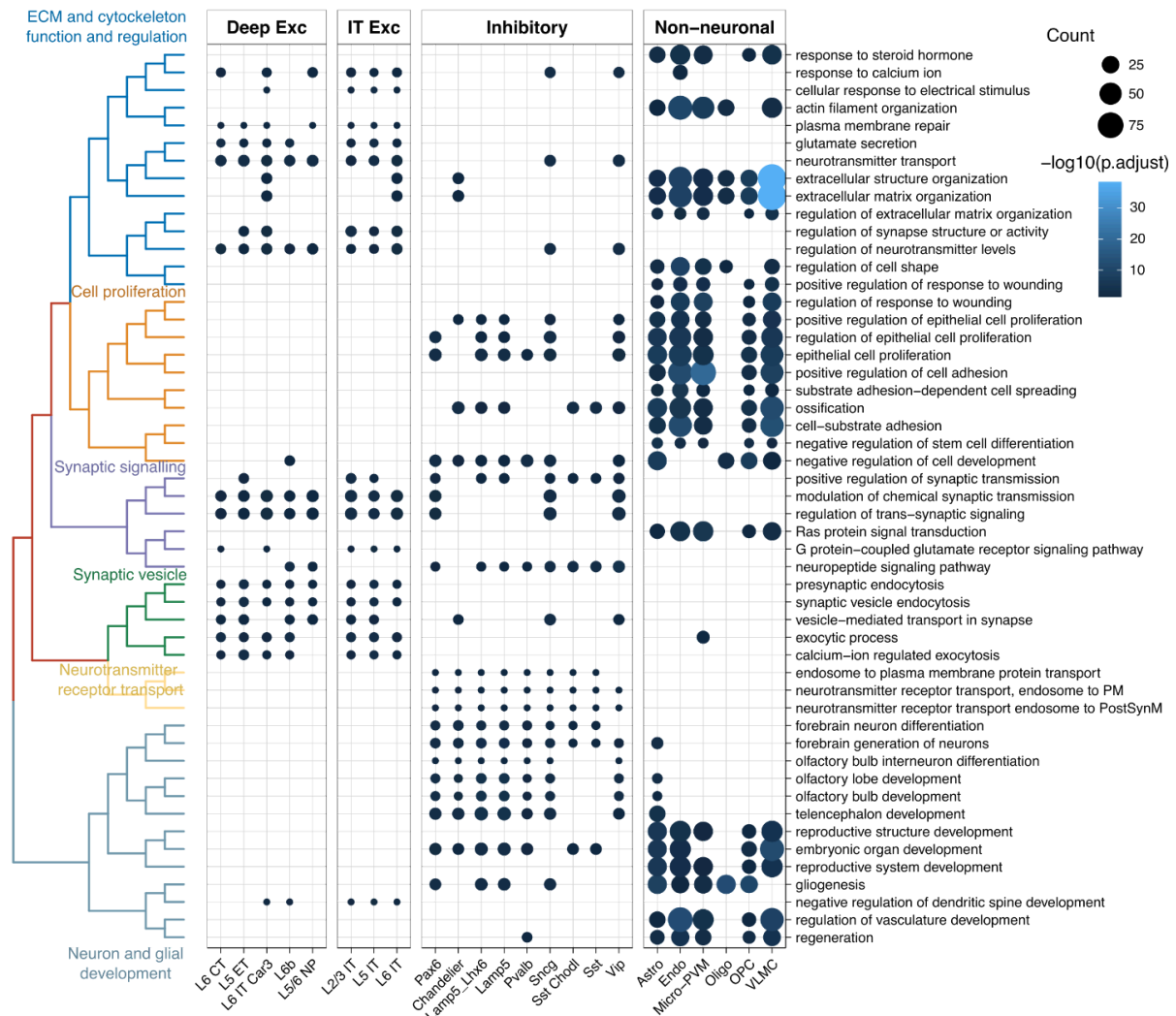
**Supplementary Figure 3 Gene specificity class conservation among O2O orthologs in the primate MTG dataset, for each species pair.** Average specificity class conservation heatmap showing the percent overlap among one-to-one orthologs in various specificity class

combinations (calculation process see methods). Pie chart showing the aggregated percentage of cell type conservation and class conservation for genes with cell type level specificity or group level specificity.



**Supplementary figure 4 Gene specificity class switching mostly happens between transcriptomically similar cell types.** (a) Showing the top cell type pairs in which cell type-specific genes switched between them. (b) Showing examples of cell type enriched genes that switched cell type between species with the top frequency. Mean expr., mean scaled expression. VLMC, vascular leptomenigeal cell; OPC, oligodendrocyte precursor cell;

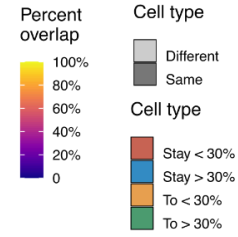
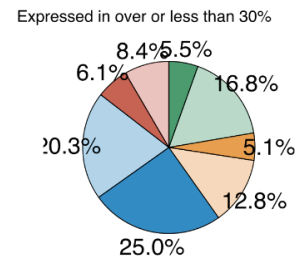
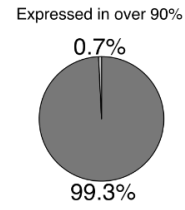
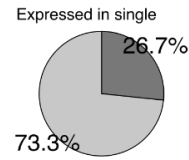
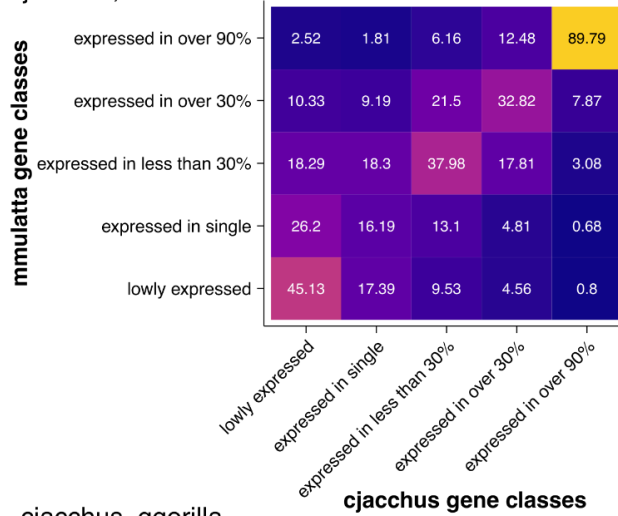
Micro-PVM, microglial cell; L6 IT, L6 intratelencephalic neurons; Endo, cerebral cortex endothelial cell; Astro, astrocyte of the cerebral cortex.



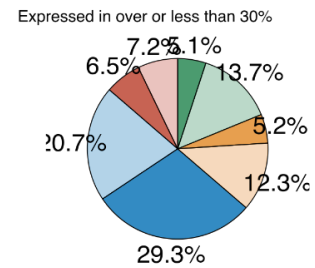
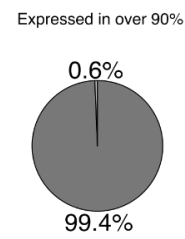
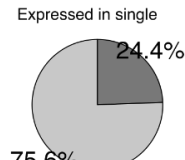
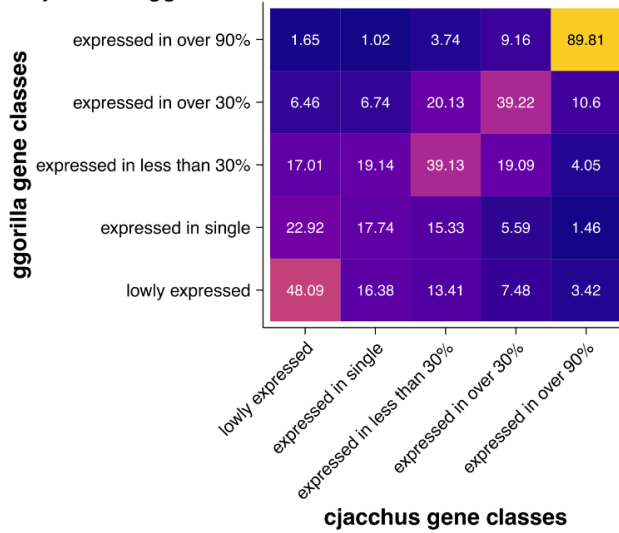
**Supplementary figure 5 GO BP enrichment of genes that stayed in the same specificity class.** Here shows the GO BP enrichment of genes that stayed enriched in the same cell type in any pair of species. Top most frequent enriched terms across each broad type were selected and plotted in the figure. Semantic similarity between the selected GO BP terms were calculated and used for hierarchical clustering, showing on the left and summarised by keywords.



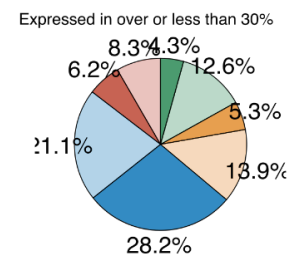
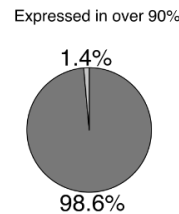
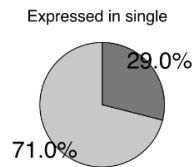
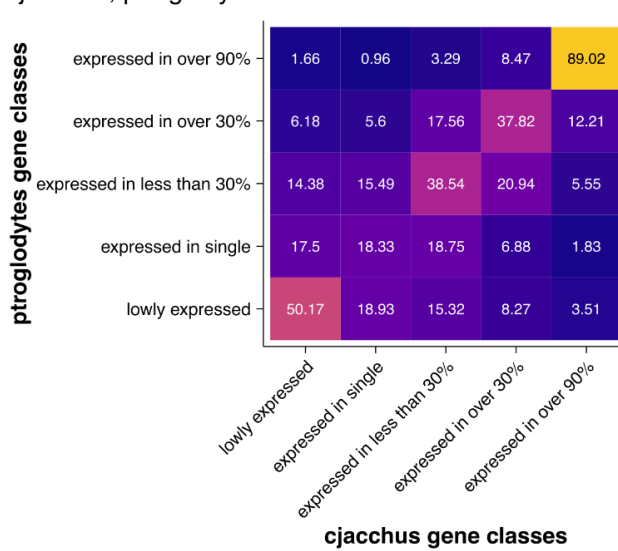
### cjacchus, mmulatta

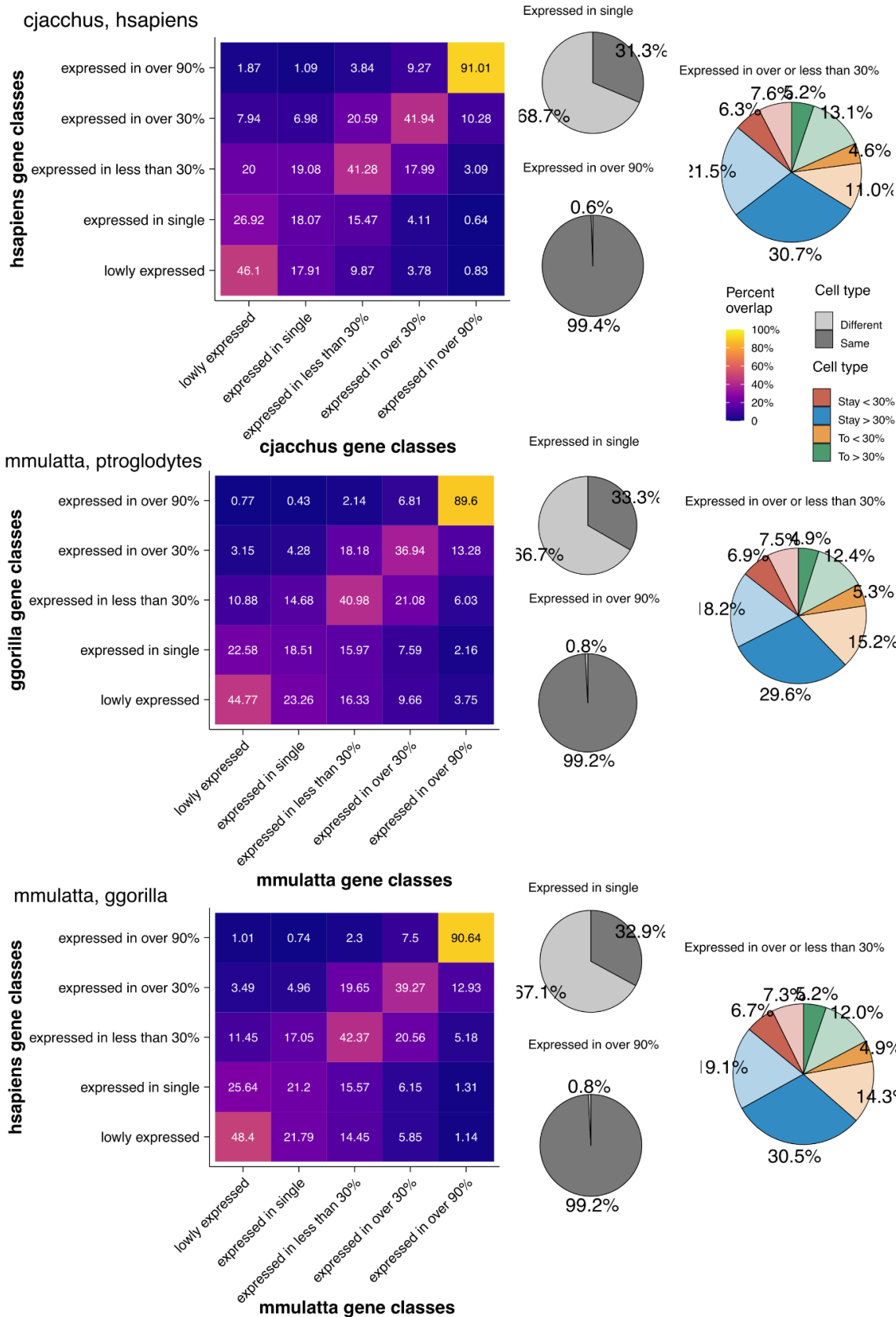


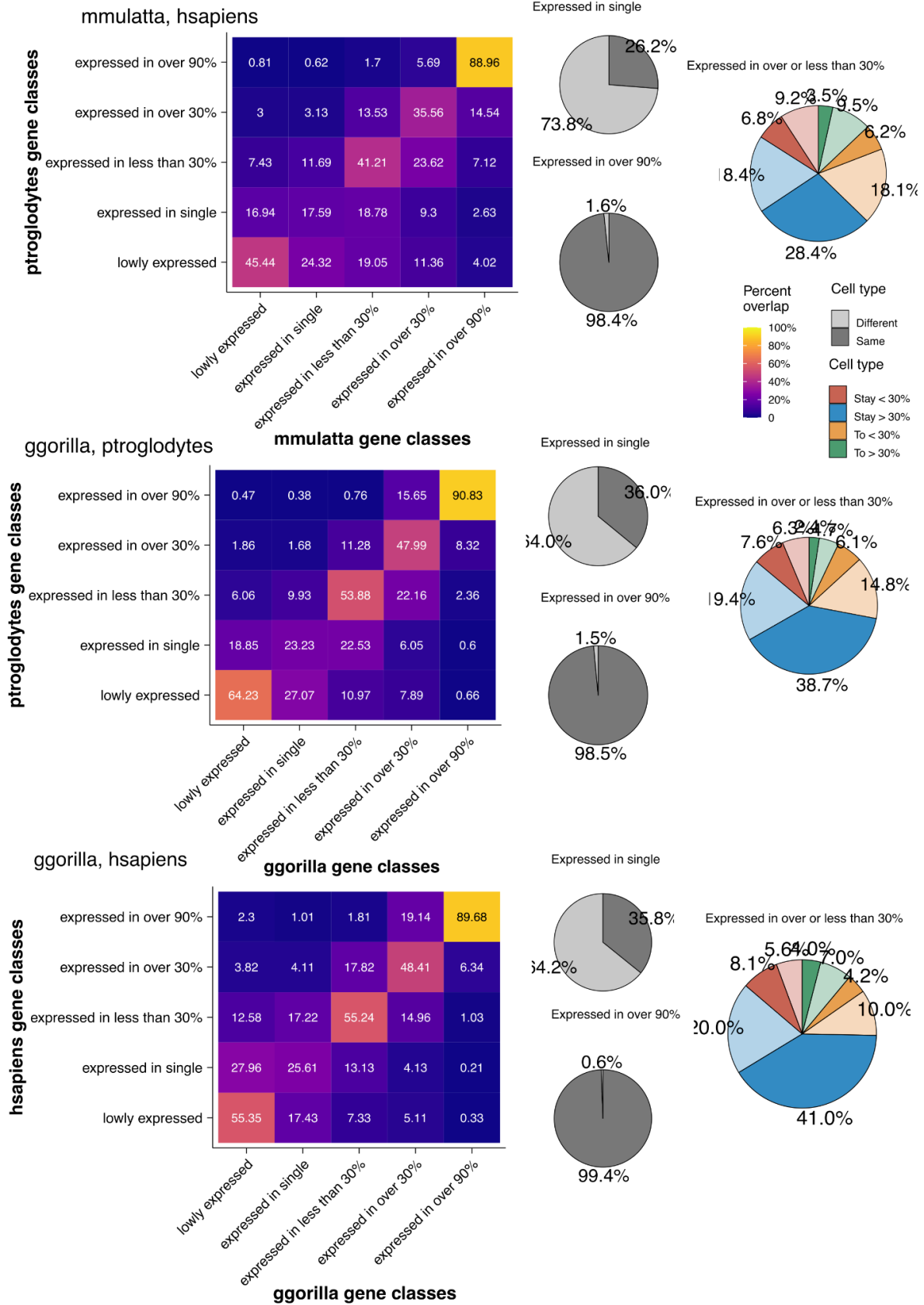
### cjacchus, ggorilla

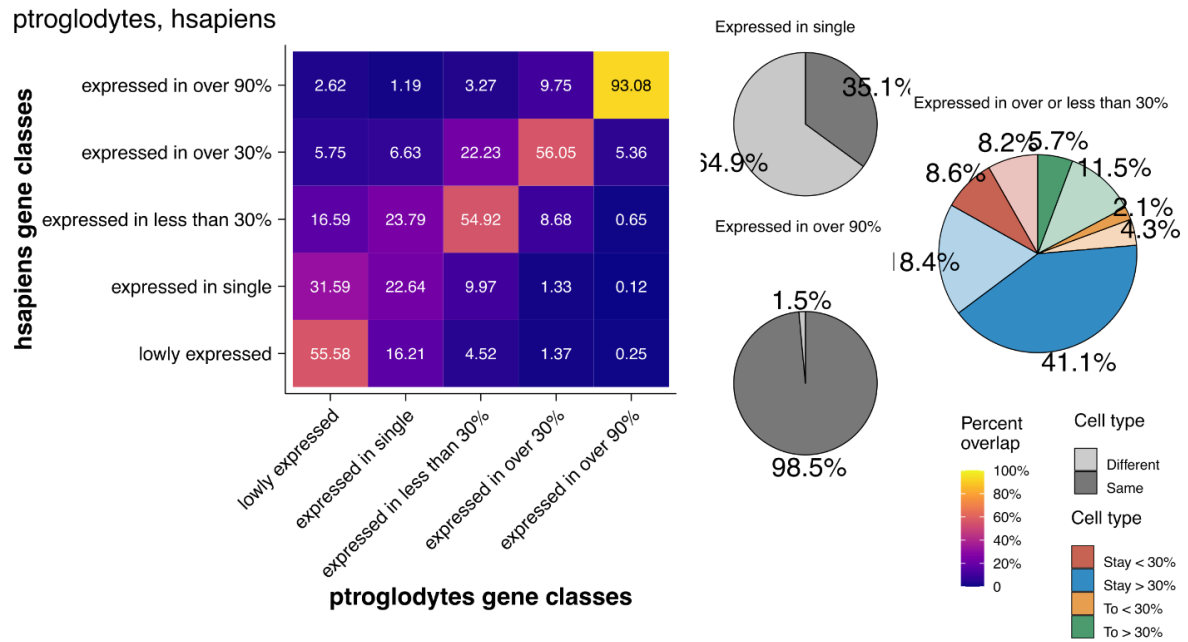


### cjacchus, ptroglodytes



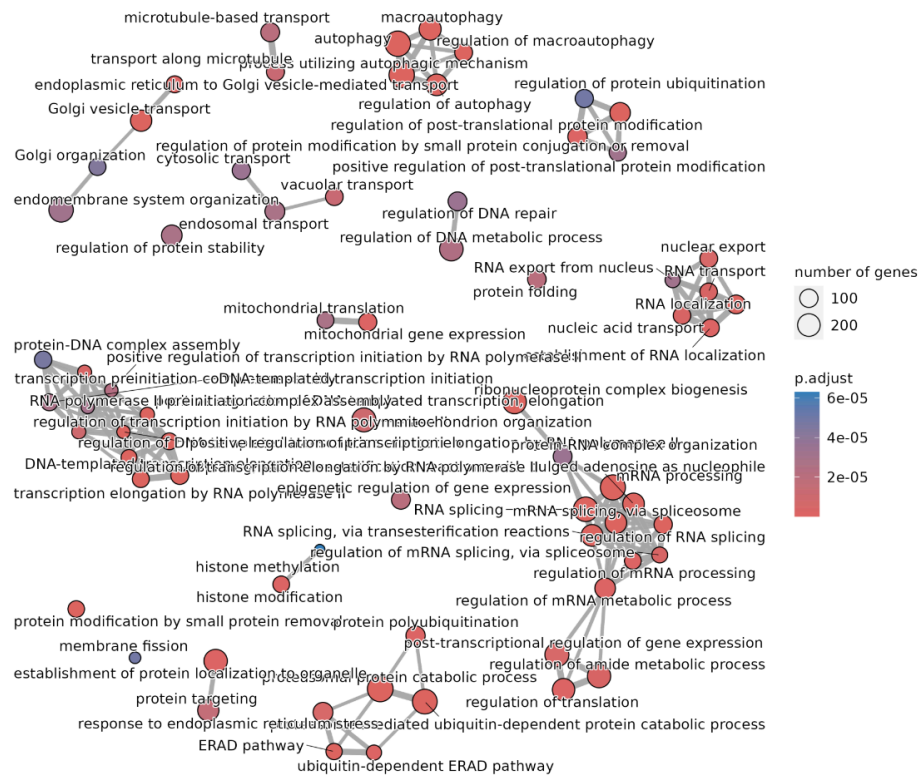




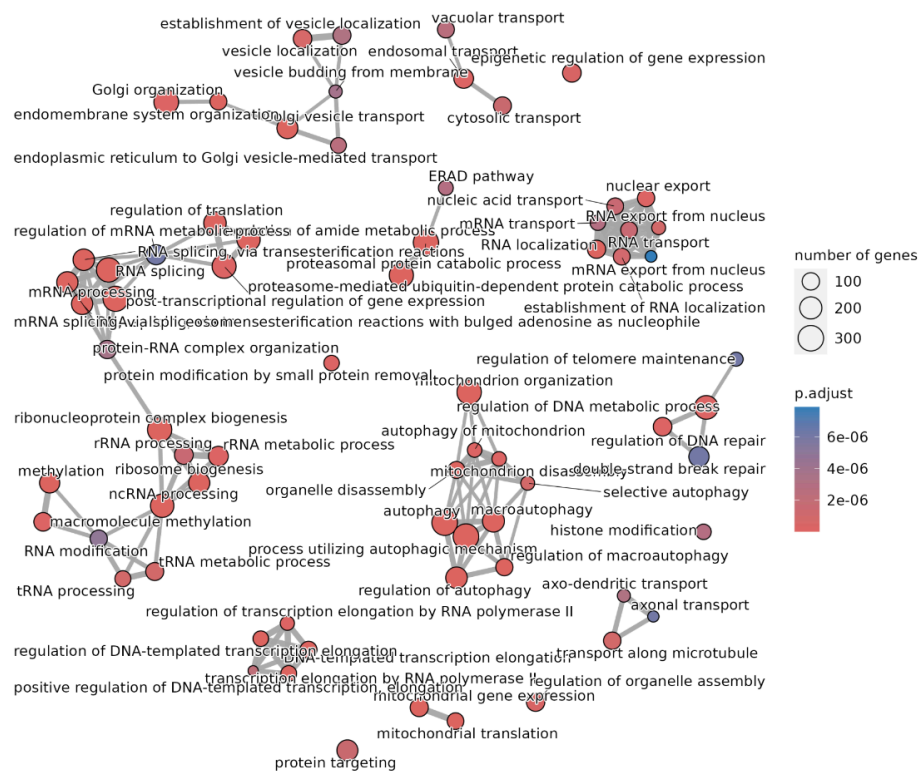


**Supplementary figure 6 Gene distribution class conservation matrix for each species pair.** Showing the average distribution class conservation heatmap (calculation process see methods). Pie chart showing the aggregated percentage of cell type conservation of genes that stayed expressed in over 90% of cell types; genes stayed expressed in a single cell type, and genes expressed in some cell types.

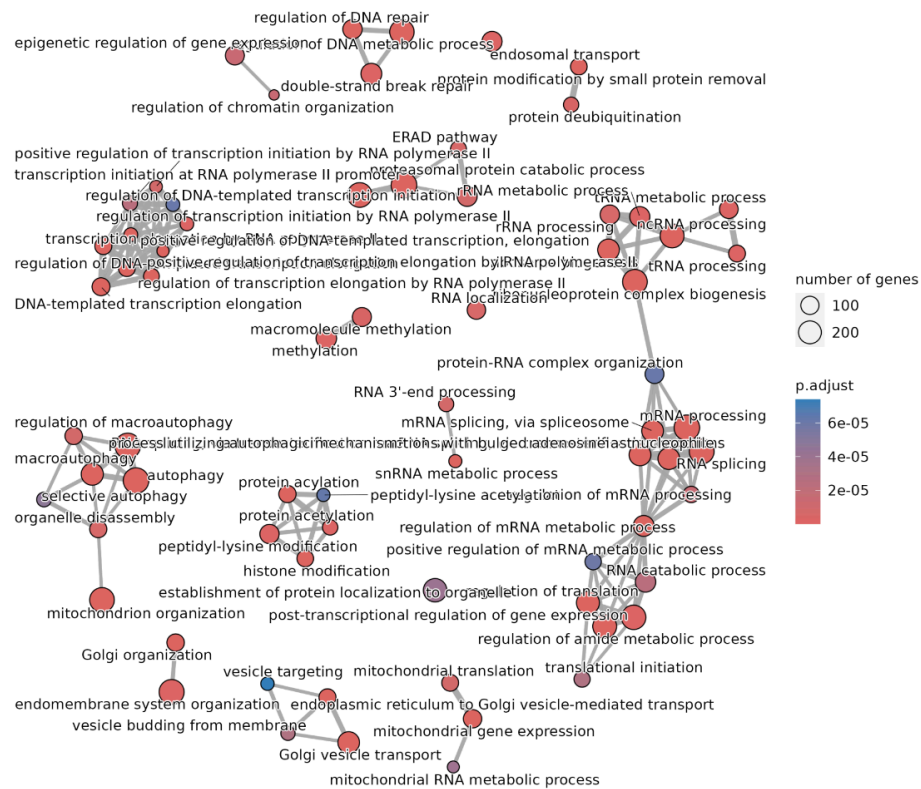
**a** *H.sapiens* and *P.troglodytes* conserved lowly specific and broadly expressed genes



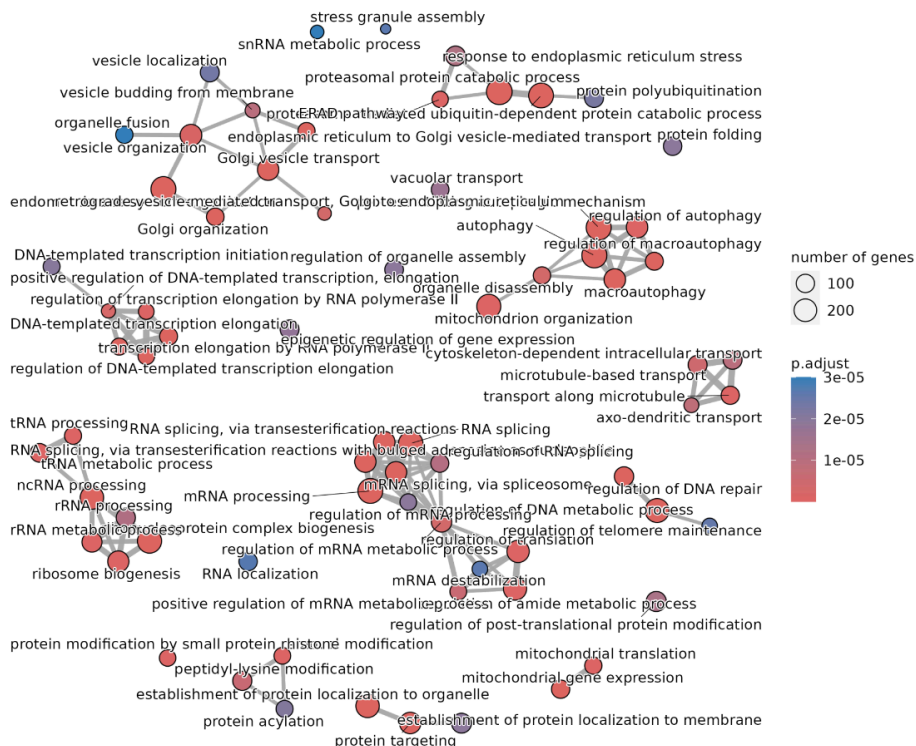
**b** *H.sapiens* and *G.gorilla* conserved lowly specific and broadly expressed genes



**c** *H.sapiens* and *M.mulatta* conserved lowly specific and broadly expressed genes



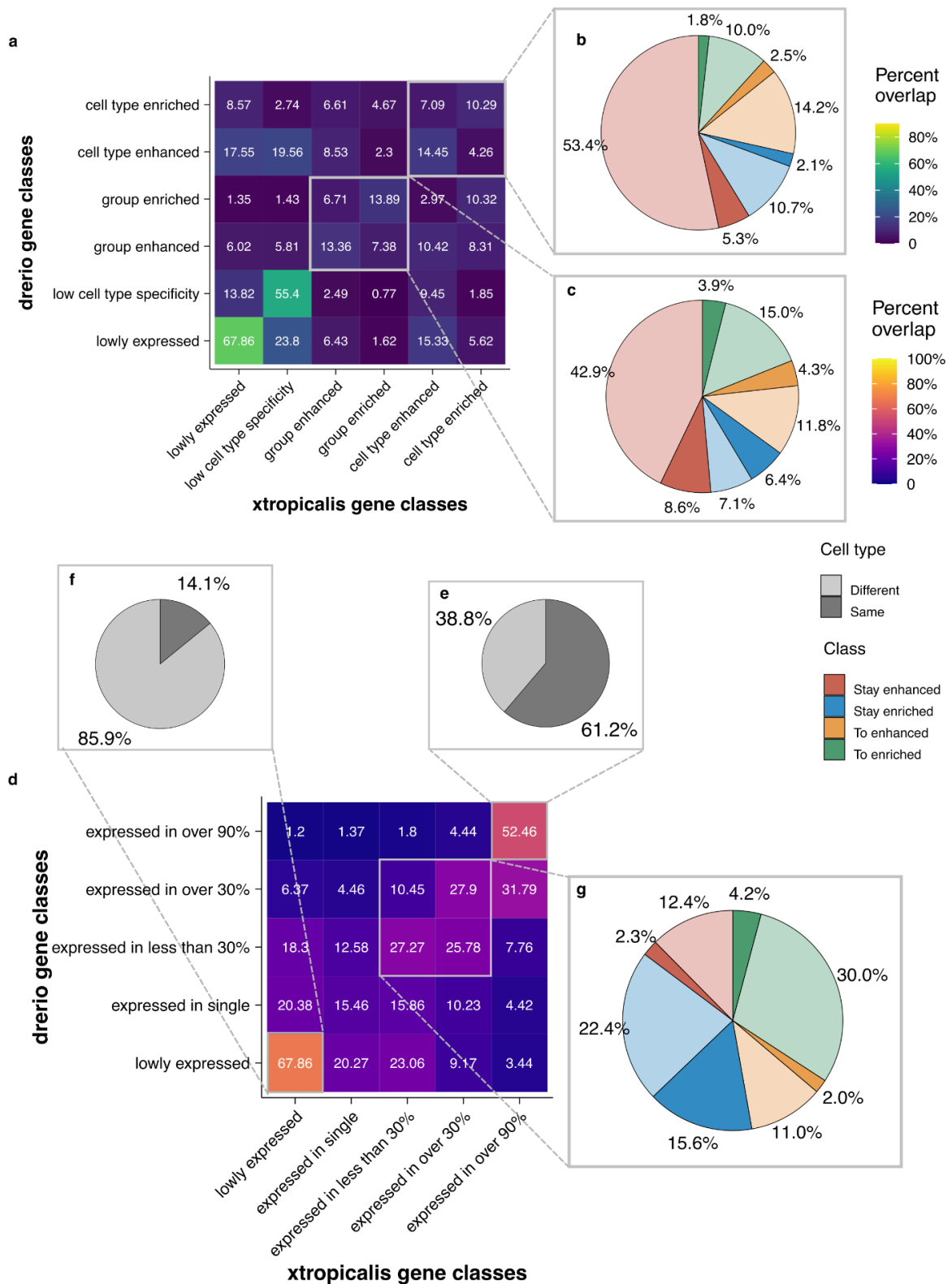
**d** *H.sapiens* and *C.jacchus* conserved lowly specific and broadly expressed genes



**Supplementary figure 7 Gene ontology enrichment results for genes stayed lowly specific and broadly expressed between human and other species in the primate MTG data. Figure made with emaplot from the R package enrichplot. Showing 70 categories. P-values were**

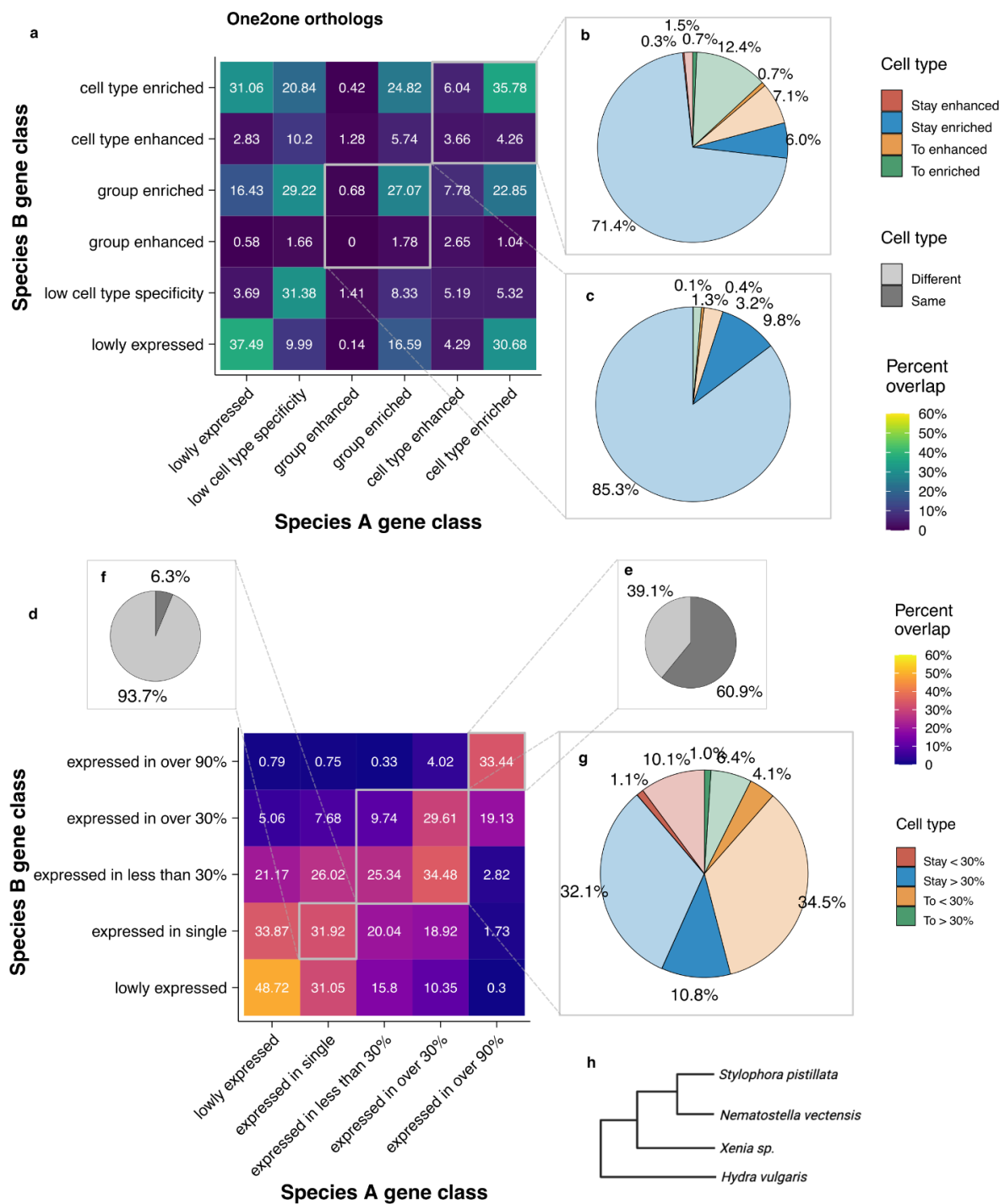


corrected for multiple comparisons using the Benjamini-Hochberg method, and results were considered significant if  $P < 0.01$ .



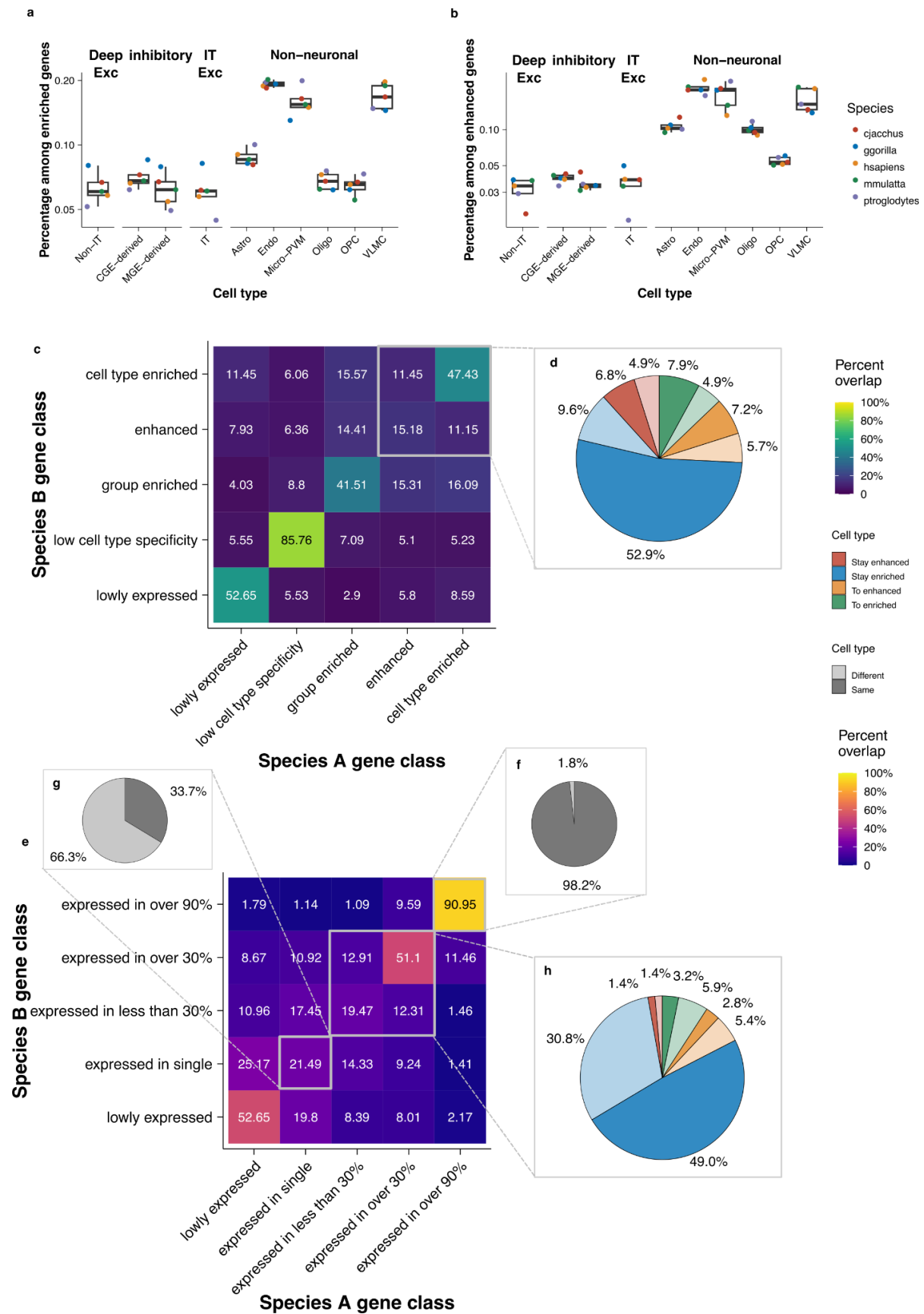
**Supplementary figure 8 Gene class conservation among O2O orthologs in the embryogenesis dataset.** (a) Average specificity class conservation heatmap showing the percent overlap among 1-2-1 orthologs in various specificity class combinations (calculation process see methods). (b)-(c) Pie chart showing the aggregated percentage of cell type conservation and class conservation for genes with (b) cell type level specificity or (c) group level specificity. (d) Average distribution class conservation heatmap (calculation process see methods). (e)-(g) Pie chart showing the aggregated percentage of cell type conservation of genes stayed expressed in (e) over 90% of cell types; genes stayed expressed in (f) a single cell type, and genes expressed in (g) some cell types.





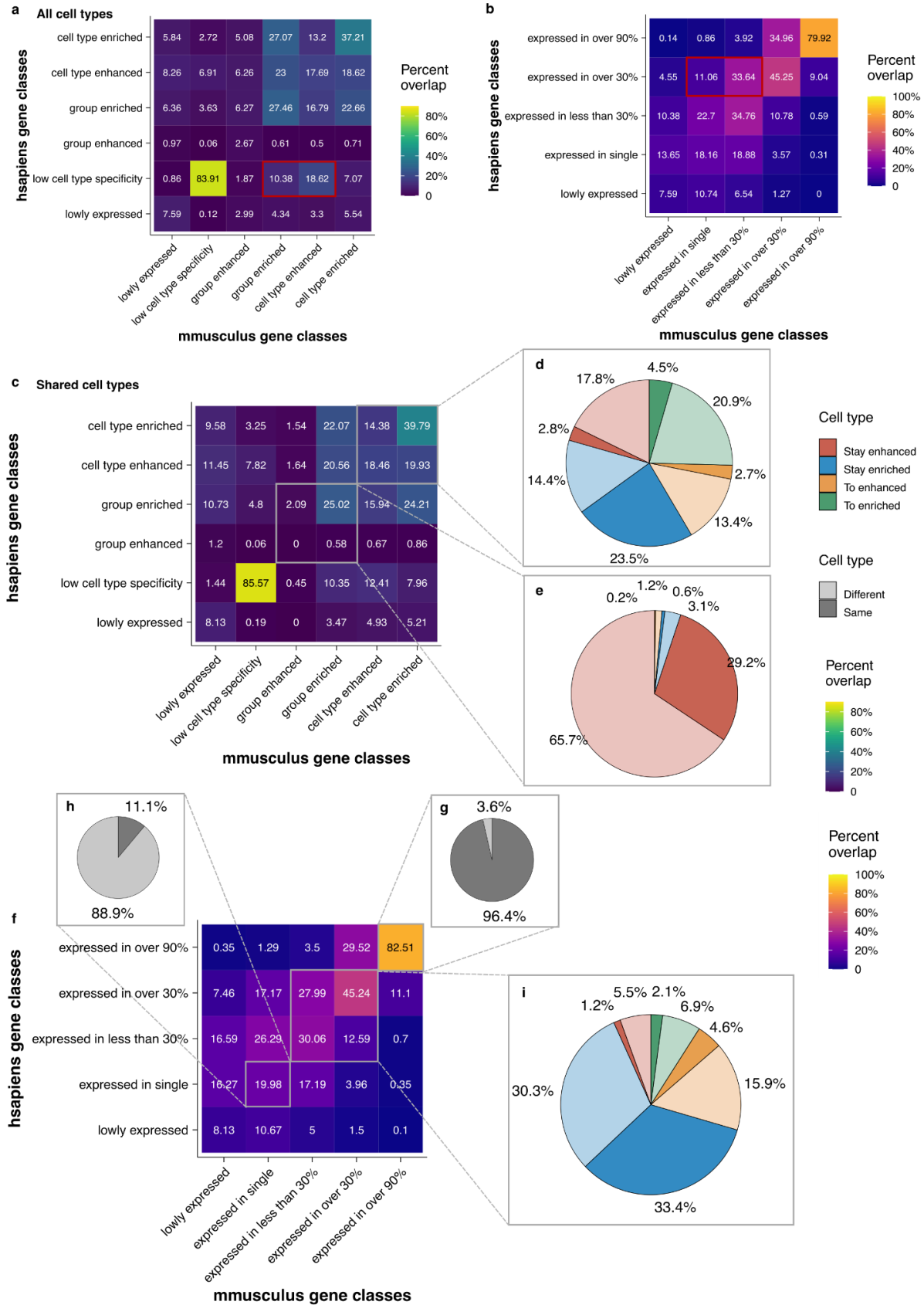
**Supplementary figure 9 Gene class conservation among O2O orthologs in the Cnidarian dataset.** (a) Average specificity class conservation heatmap showing the percent overlap among 1-2-1 orthologs in various specificity class combinations (calculation process see

methods). (b)-(c) Pie chart showing the aggregated percentage of cell type conservation and class conservation for genes with (b) cell type level specificity or (c) group level specificity. (d) Average distribution class conservation heatmap (calculation process see methods). (e)-(g) Pie chart showing the aggregated percentage of cell type conservation of genes stayed expressed in (e) over 90% of cell types; genes stayed expressed in (f) a single cell type, and genes expressed in (g) some cell types. (h) species tree of this dataset.



**Supplementary figure 10 Gene class conservation among O20 orthologs in primate MTG dataset using reduced cell type granularity.** Neurons are grouped by reduced granularity but

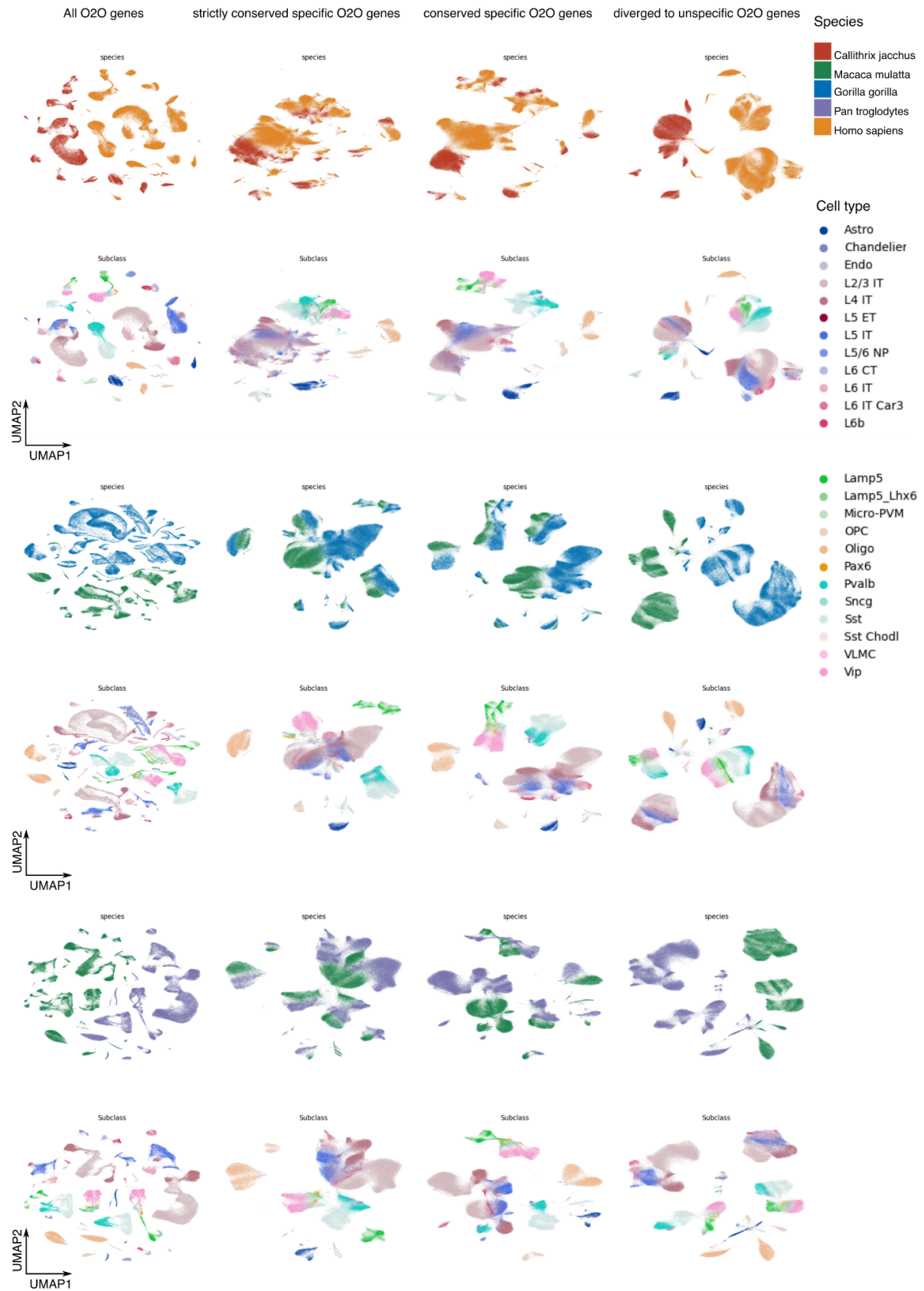
not non-neuronal cells. (a)-(b) Number of cell type-specific genes across species for each cell type. (c) Average specificity class conservation heatmap showing the percent overlap among one-to-one orthologs in various specificity class combinations (calculation process see methods). (d) Pie chart showing the aggregated percentage of cell type conservation and class conservation for genes with cell type level specificity. (e) Average distribution class conservation heatmap (calculation process see methods). (f)-(h) Pie chart showing the aggregated percentage of cell type conservation of genes stayed expressed in (f) over 90\% of cell types; genes stayed expressed in (g) a single cell type and genes expressed in (h) some cell types.



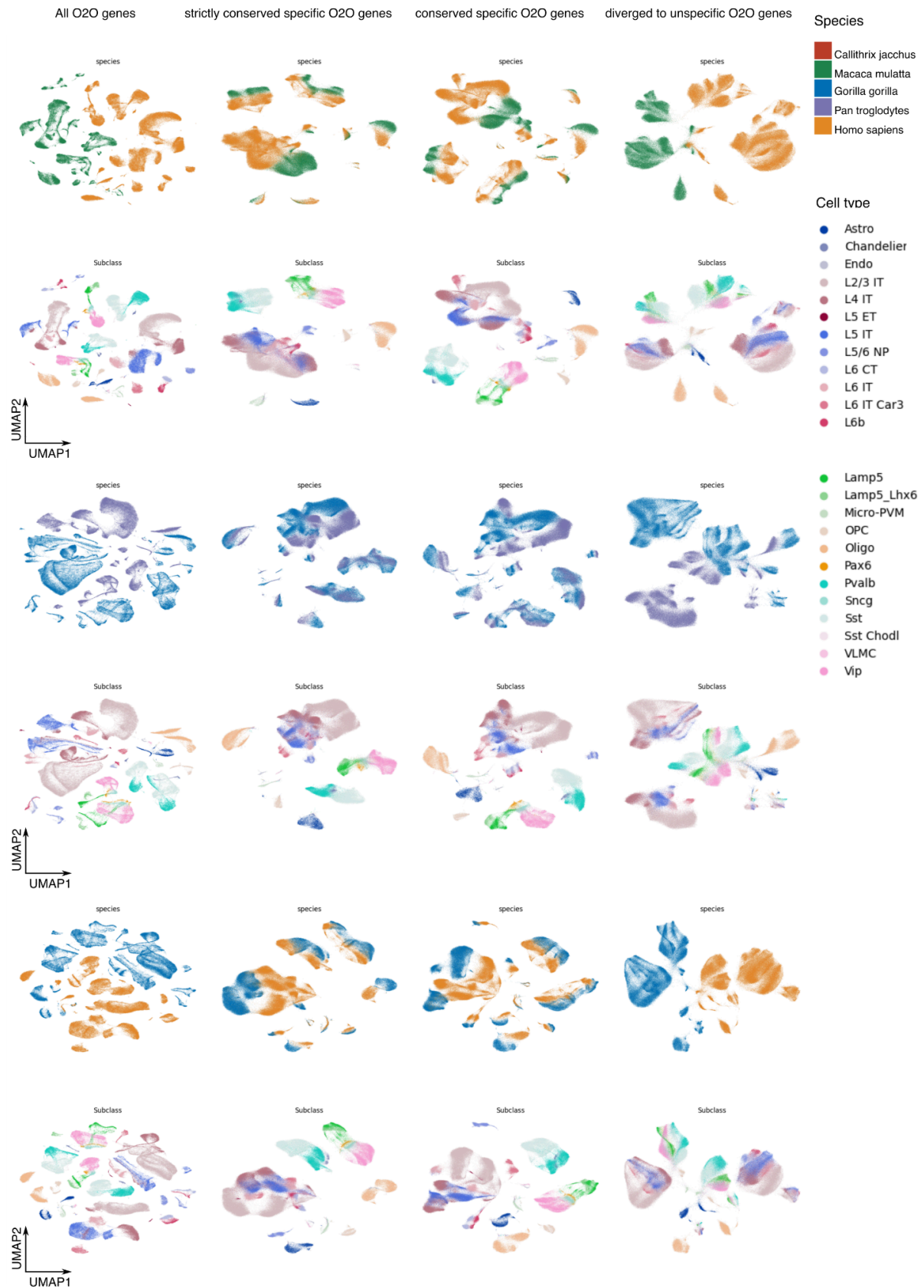
**Supplementary figure 11 Gene class conservation among O2O orthologs in the human-mouse bone marrow dataset using all cell types, compared with using only shared**

**cell types.** (a) Specificity class conservation heatmap showing the percent overlap among 1-2-1 orthologs in various specificity class combinations using all cell types (b) distribution class conservation heatmap using all cell types. Red boxes indicate the stripe pattern, indicating genes with specificity in mouse but not human, attributed to mouse data-specific cell types. (c) Specificity class conservation heatmap showing the percent overlap among 1-2-1 orthologs in various specificity class combinations using shared cell types. (d)-(e) Pie chart showing the aggregated percentage of cell type conservation and class conservation for genes with cell type level specificity (d) or group level specificity (e). (f) Distribution class conservation heatmap using shared cell types. (g)-(i) Pie chart showing the aggregated percentage of cell type conservation of genes stayed expressed in over 90% of cell types (g); genes stayed expressed in a single cell type (h,) and genes expressed in some cell types (i).



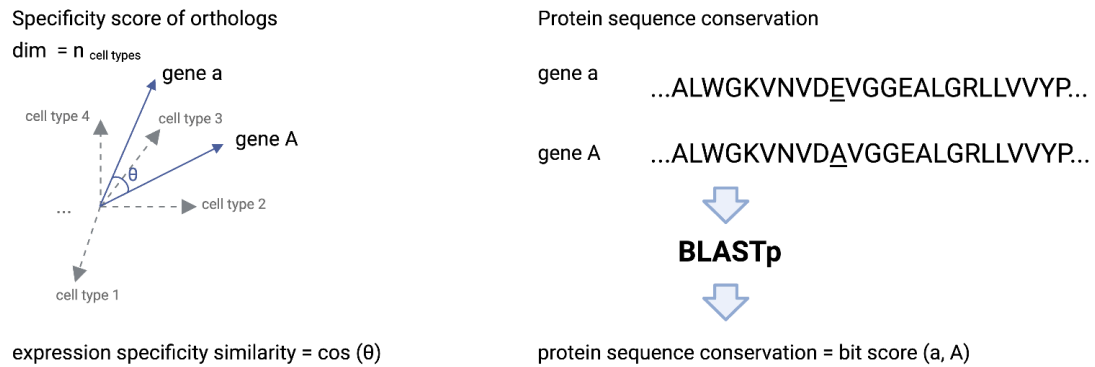






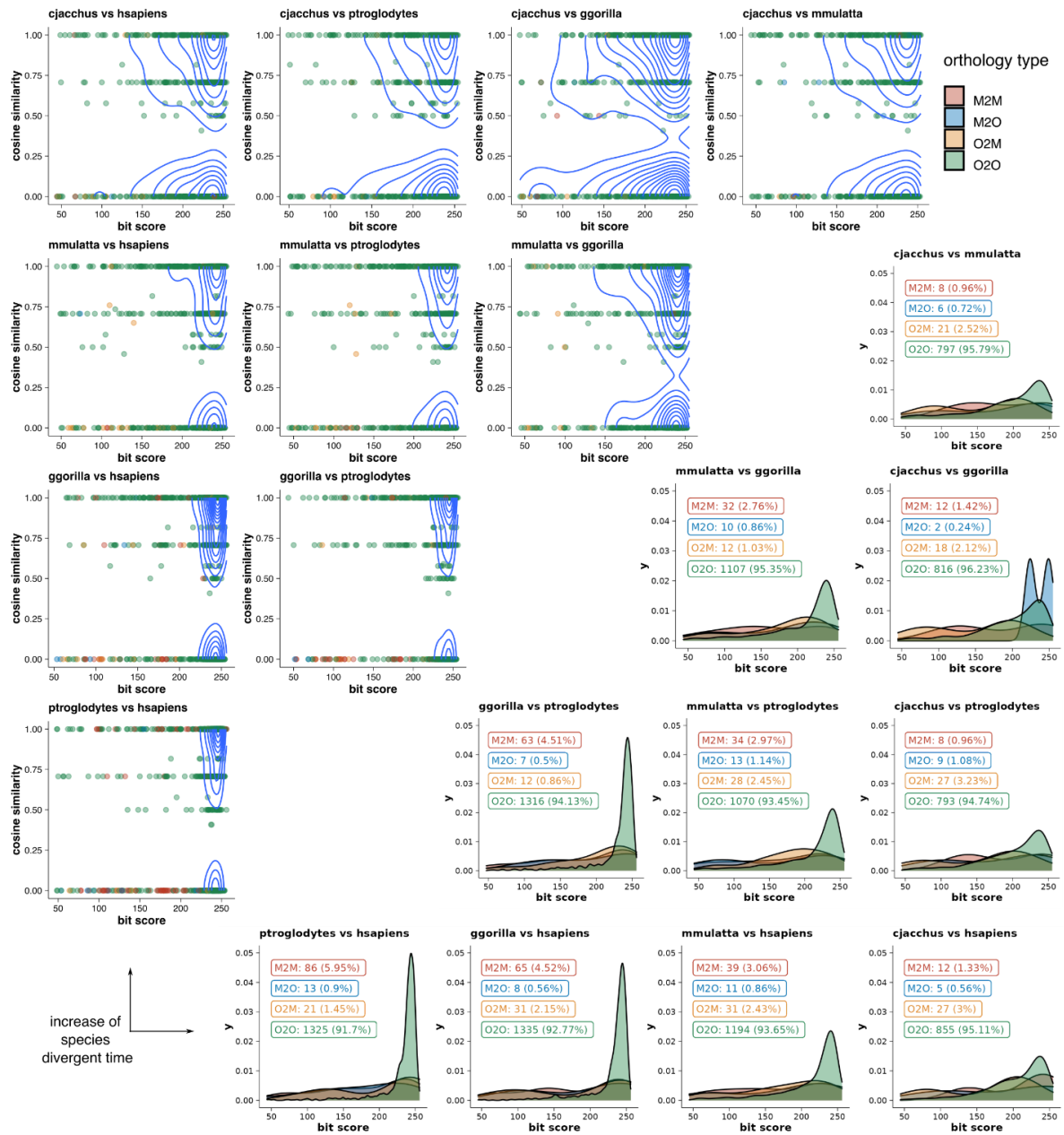
**Supplementary figure 12 Gene class switching explains species effect on transcriptomic space for each species pair. (a) UMAP visualisation using different sets of O2O orthologs**

for all species data. Showing the strong species effect using all genes, the reduced species effect using (strictly) conserved specific genes, as well as a purified species effect observed in data using diverged to unspecific genes. O2O, one-to-one; UMAP, Uniform Manifold Approximation and Projection; PC, principal component.



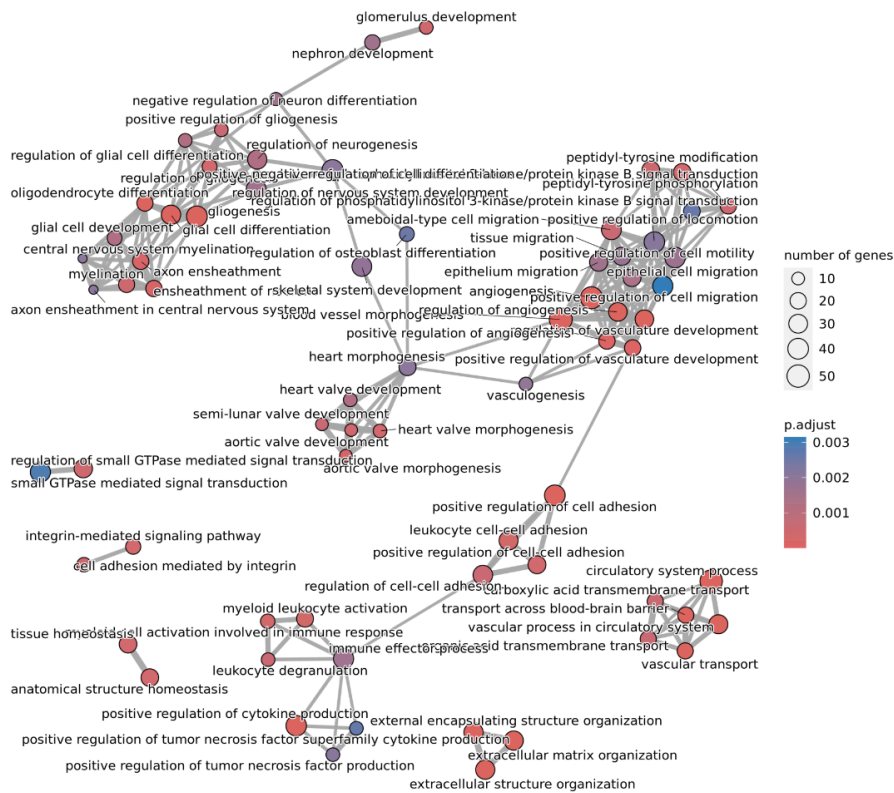
**Supplementary figure 13 Schematic of the comparison of expression specificity similarity and sequence conservation of orthologs.** Gene A and gene a represent orthologs from different species. Expression specificity similarity is calculated with the cosine similarity of gene specificity scores across cell types. Protein sequence conservation is the bit score of BLASTp.

### Cell type or group enhanced genes

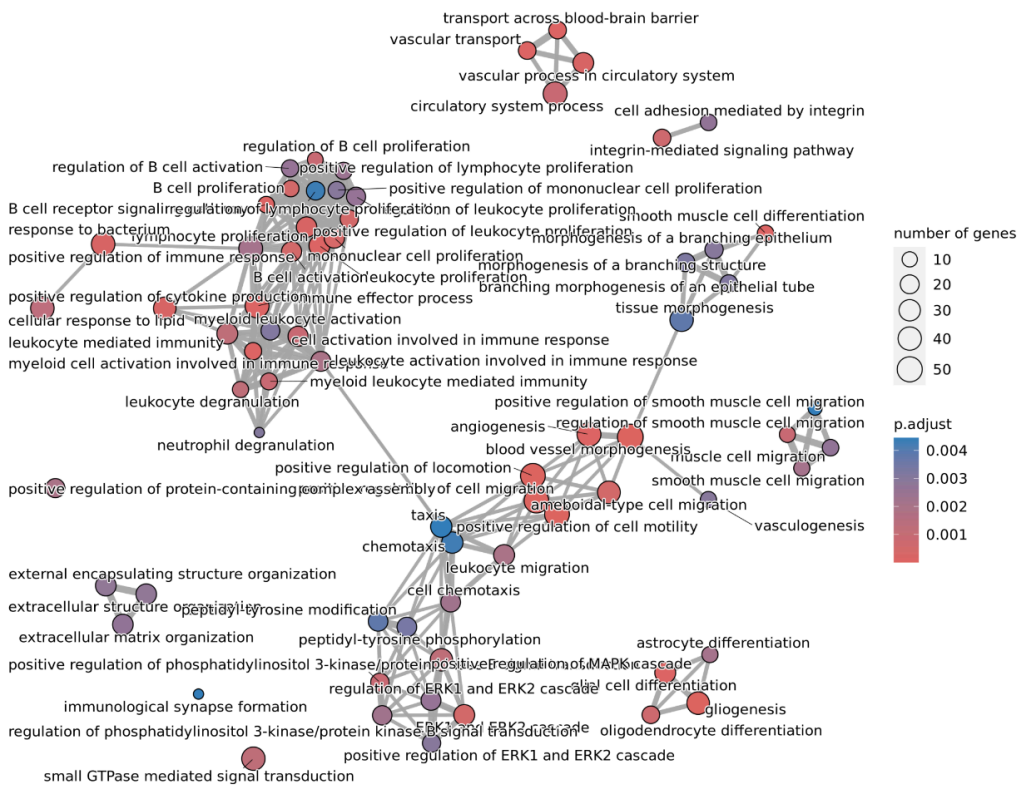


**Supplementary figure 14 Comparing sequence conservation with cell type expression specificity similarity across species distance for enhanced genes.** The top left triangle plots the sequence conservation (bit score) against the expression similarity (cosine similarity) for each pair of cell type or group enhanced orthologs in each pair of species. The bottom right triangle shows the distribution of bit scores for different types of orthologs for the same set of genes. O2O: one-to-one, O2M: one-to-many, M2O: many-to-one, M2M: many-to-many.

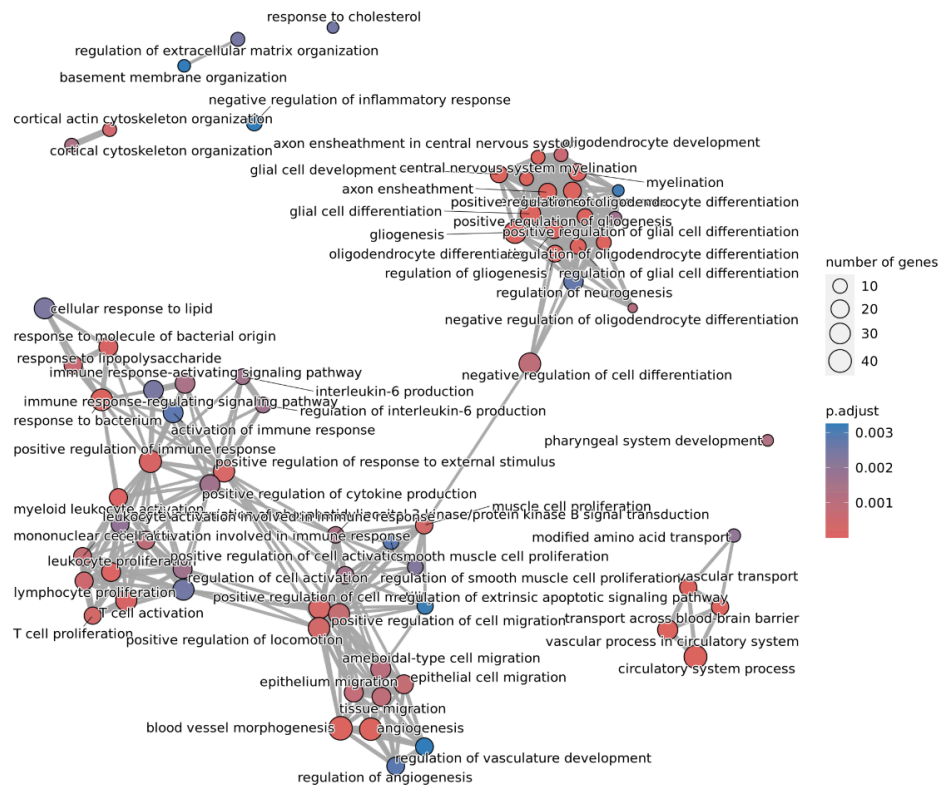
**a** *H.sapiens* and *P.troglodytes* conserved cell type-specific and high sequence conservation genes



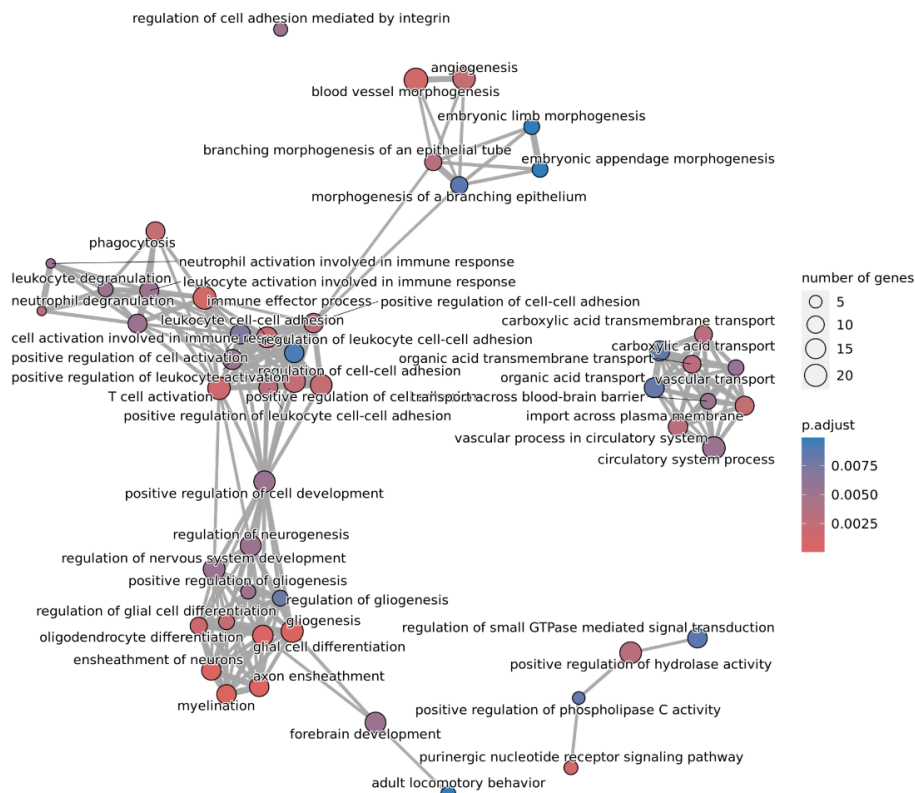
**b** *H.sapiens* and *G.gorilla* conserved cell type-specific and high sequence conservation genes



**c** *H.sapiens* and *M.mulatta* conserved cell type-specific and high sequence conservation genes



**d** *H.sapiens* and *C.jacchus* conserved cell type-specific and high sequence conservation genes



**Supplementary figure 15 Gene ontology enrichment results for genes that have complete specificity conservation and high sequence similarity between human and other species.**



**a** *D. rerio* vs *X. tropicalis* cell type or group enriched orthologs

Orthology type

- M-2-M
- M-2-1
- 1-2-M
- 1-2-1

**b**

Orthology type

- M-2-M
- M-2-1
- 1-2-M
- 1-2-1

**c** *D. rerio* vs *X. tropicalis* cell type or group enhanced orthologs

Orthology type

- M-2-M
- M-2-1
- 1-2-M
- 1-2-1

**d**

Orthology type

- M-2-M
- M-2-1
- 1-2-M
- 1-2-1

Orthology type	Count	Percentage
M-2-M	16	23.53%
M-2-1	4	5.88%
1-2-M	5	7.35%
1-2-1	43	63.24%

Orthology type	Count	Percentage
M-2-M	15	5.98%
M-2-1	24	9.56%
1-2-M	2	0.8%
1-2-1	210	83.67%

