

1 **Cross-species cell-type assignment of single-cell RNA-seq**
2 **by a heterogeneous graph neural network**

3
4
5 Xingyan Liu^{1,2†}, Qunlun Shen^{1,2†} and Shihua Zhang^{1,2,3,4*}

6 ¹NCMIS, CEMS, RCSDS, Academy of Mathematics and Systems Science,
7 Chinese Academy of Sciences, Beijing 100190, China;

8 ²School of Mathematical Sciences, University of Chinese Academy of Sciences,
9 Beijing 100049, China;

10 ³Center for Excellence in Animal Evolution and Genetics, Chinese Academy of
11 Sciences, Kunming 650223, China;

12 ⁴Key Laboratory of Systems Biology, Hangzhou Institute for Advanced Study,
13 University of Chinese Academy of Sciences, Chinese Academy of Sciences,
14 Hangzhou 310024, China.

15 †These authors contributed equally to this work.

16 *To whom correspondence should be addressed. Tel/Fax: +86 01 82541360;
17 Email: zsh@amss.ac.cn.

18

19 **Abstract**

20 Cross-species comparative analyses of single-cell RNA sequencing (scRNA-
21 seq) data allow us to explore, at single-cell resolution, the origins of cellular
22 diversity and the evolutionary mechanisms that shape cellular form and function.
23 Here, we aimed to utilize a heterogeneous graph neural network to learn
24 aligned and interpretable cell and gene embeddings for cross-species cell type
25 assignment and gene module extraction (CAME) from scRNA-seq data. A
26 systematic evaluation study on 649 pairs of cross-species datasets showed that
27 CAME outperformed six benchmarking methods in terms of cell-type
28 assignment and model robustness to insufficiency and inconsistency of
29 sequencing depths. Comparative analyses of the major types of human and
30 mouse brains by CAME revealed shared cell type-specific functions in
31 homologous gene modules. Alignment of the trajectories of human and
32 macaque spermatogenesis by CAME revealed conservative gene expression
33 dynamics during spermatogenesis between humans and macaques. Owing to
34 the utilization of non-one-to-one homologous gene mappings, CAME made a
35 significant improvement on cell-type characterization cross zebrafish and other
36 species. Overall, CAME can not only make an effective cross-species
37 assignment of cell types on scRNA-seq data but also reveal evolutionary
38 conservative and divergent features between species.

39

40 **Key words**

41 Cross-species; cell-type assignment; gene module extraction; single-cell RNA
42 sequencing; a heterogeneous graph neural network

43

44 Introduction

45 Single-cell RNA sequencing (scRNA-seq) has rapidly emerged as a powerful
46 tool to characterize a large number of single-cell transcriptomes in different
47 tissues, organs, and species [2]. It not only deepens our knowledge of cells but
48 also provides novel insights into evolutionary and developmental biology [3].
49 Cross-species integration and comparison of scRNA-seq datasets allow us to
50 explore, at single-cell resolution, the origins of cellular diversity and the
51 evolutionary mechanisms that shape cellular form and function [3-11].

52 Cell-type assignment (or cell typing) and data integration are both vital steps
53 involved in these analyses. For the cell-type assignment, a traditional approach
54 includes three steps: clustering single-cells, performing differentially expression
55 analysis to find cluster-specific genes, and matching these genes with known
56 markers. However, this strategy fails when different cell types are clustered into
57 one group, and when analyzing many non-model species that lack prior
58 knowledge of cell-type biomarkers. Several tools have been developed for this
59 task recently. Some existing approaches like CellAssign [12] and scCATCH [13]
60 require prior knowledge of cell type-specific markers. Some like SingleCellNet
61 [14] and SciBet [15] were designed based on a reference dataset and can
62 achieve the cell-type assignment without providing marker information. Besides,
63 several methods designed for data integration can also achieve cell-type
64 assignment by transferring labels from the reference dataset. Seurat-v3 [16]
65 combines canonical correlation analysis and mutual nearest neighbors to
66 perform data integration and label transfer based on ‘anchors’. Cell BLAST [17]
67 and ItClust [18] make use of deep neural networks for both cell-type querying
68 and cell embedding. LIGER [19] and CSMF [20] extract the common and
69 private features of two datasets respectively by joint non-negative matrix
70 factorization to achieve cell alignment across datasets and omics.

71 Despite all the progress, a tool for effective and robust cross-species
72 integration and comparison is still immature and in demand. There are several
73 computational challenges to be overcome. First, it is hard to determine cell
74 identities for non-model species that lack prior knowledge of cell-type
75 biomarkers, and most of the methods may fail when generalizing to cross-
76 species label transfer. Second, many biological and technical factors, such as
77 transcriptome variation between species, different experimental protocols, and
78 inconsistent sequencing depths, can make cross-species data integration and
79 comparison even more difficult. Third, homologous cell-type alignment requires
80 quantifying the similarities of gene expression profiles, which usually vary
81 across distinct normalizations and gene selections [3]. Fourth, cross-species
82 cellular alignment is usually based on homologous genes and current
83 approaches are mostly restricted to one-to-one homologies shared by both
84 organisms [3, 5-11], where non-one-to-one homologous genes characterizing
85 cell-type conservative features could be lost. Lastly, evolutionary divergences

86 are thought to be caused by transcriptional changes of groups of genes that
87 evolve in a modular fashion and are controlled by transcription factors [21].
88 Extraction and comparison of gene modules between species will provide deep
89 insights into evolutionary conservation and divergences [11, 22, 23].

90 To this end, we developed a heterogeneous graph neural network model to
91 achieve the aligned and interpretable cell and gene embeddings for cross-
92 species cell-type assignment and gene module extraction (CAME). A
93 systematic evaluation study on 649 pairs of cross-species datasets showed that
94 CAME outperformed six benchmarking methods in terms of cell-type prediction,
95 and model robustness to insufficiency and inconsistency of sequencing depths.
96 Comparative analyses of the major types of human and mouse brains by CAME
97 revealed shared cell type-specific functions in homologous gene modules. An
98 alignment of the trajectories of human and macaque spermatogenesis by
99 CAME revealed the conservative gene expression dynamics during
100 spermatogenesis between humans and macaques. Owing to the utilization of
101 non-one-to-one homologous gene mappings, CAME made a significant
102 improvement on cell-type characterization across long-distant species. Overall,
103 CAME can not only make an effective cross-species assignment of cell types
104 on scRNA-seq data but also reveal evolutionary conservative and divergent
105 features between species.

106 Results

107 Overview of CAME

108 CAME takes two scRNA-seq datasets from different species, along with their
109 homologous gene mappings as input. One dataset with cell-type labels is taken
110 as the reference and the other whose cell types need to be assigned is the
111 query (**Figure 1A**). CAME encodes these two expression matrices and the
112 mappings of homologous genes as a heterogeneous graph, where each node
113 acts as either a cell or a gene, while a cell-gene edge indicates a non-zero
114 expression of the gene in that cell, and an edge between a pair of genes
115 indicates the homology between each other. Note that one-to-many and many-
116 to-many homologies are allowed as well. Besides, CAME adopts single-cell
117 networks pre-computed from reference and query datasets using the k-nearest-
118 neighbor (KNN) method, respectively, where a cell-cell edge indicates this pair
119 of cells have similar transcriptomes with each other (**Methods**).

120 CAME adopts a heterogeneous graph neural network to embed each node
121 into a low-dimensional space (**Methods, Figure 1B**). For the initial cell
122 embeddings, CAME takes the expression profiles followed by linear
123 transformation with a non-linear activation function. While for the initial gene-
124 embeddings, CAME aggregates the expression profiles (called “message”)
125 from its neighbor cells which expressed it, and then treats them with linear
126 transformation and non-linear activation, as done for cells (**Methods**). Then the

127 initial embeddings are input to two parameter-sharing graph convolution layers
128 with heterogeneous edges and nodes. As a result, cells with more co-expressed
129 genes are more likely to exchange the embedding message with each other,
130 thus be encoded with similar embeddings; the same principle applies to genes.
131 CAME further employs a heterogeneous graph attention mechanism [25] to
132 classify cells with embeddings of their neighbor genes as input, where each cell
133 pays a distinct level of attention to each certain neighbor gene (**Methods**,
134 **Figure 1B**). High attention paid by a cell to a gene implies that the gene is of
135 relatively much importance for the cell to be characterized.

136 We note that a reference cell could be assigned with multiple labels in
137 different hierarchies, and a cell type in query species might correspond to
138 multiple ones in the reference. Thus, multi-label classification can be helpful to
139 depict the state of a cell. CAME calculates the cross-entropy between the
140 predicted cell-type probabilities and the true labels for the reference data to
141 obtain both the multi-class and the multi-label loss, and sums them up as the
142 training loss. Finally, CAME minimizes it by the backpropagation algorithm
143 (**Methods**). The training process of CAME is semi-supervised in an end-to-end
144 manner. We found that the training process was quite stable, and the model
145 tended to be well trained before 200-300 epochs (**Supplementary Figure S1A**).
146 Besides, CAME introduces the adjusted mutual information (AMI) between the
147 predicted labels and pre-clustered ones of query cells to automatically
148 determine the model checkpoint for downstream analysis (**Methods** and
149 **Supplementary Figure S1A**). Ablation experiments demonstrated that six key
150 factors adopted by CAME play roles in improving the prediction performance
151 (**Supplementary Figure S1B**).

152 CAME outputs the quantitative cell-type assignment for each query cell, that
153 is, the probabilities of cell types that exist in the reference species, which
154 enables the identification of the unresolved cell states in the query data. For
155 most cells with homologous cell types in the reference, CAME assigns them
156 with a maximal probability approximating 1. While for those unobserved cell
157 types or states, CAME would assign them to their analogs with relatively low
158 confidences (**Supplementary Figure S2**). Besides, CAME gives the aligned
159 cell and gene embeddings across species, which facilitates low-dimensional
160 visualization and joint gene module extraction (**Methods, Figure 1D**).

161 **CAME showed superior accuracy and robustness for cell-type** 162 **assignment compared to benchmarking methods**

163 We collected 54 scRNA-seq datasets from five tissues across seven different
164 species including human, macaque, mouse, chick, turtle, lizard, and zebrafish
165 (**Methods, Supplementary Figure S3A and Supplementary Table S1**) and
166 found that more than a half of the homologous genes between zebrafish and
167 other species are not one-to-one matched (**Supplementary Figure S3B**).
168 Besides, the proportion of non-one-to-one homologies between highly
169 informative gene (HIG) sets with one associated with zebrafish [26] was

170 significantly higher than that of other cross-species dataset pairs (60%-75%
171 versus 15%-40%, **Supplementary Figure S3C**). And ablation study shows that,
172 when excluding non-one-to-one homologies, the cell-typing accuracy of CAME
173 suffered a significant drop (ranging from 1.5% to 8.7% for different species-
174 pairs, 6.26% on average, with p-value = $7.8e-23$) on the zebrafish-associated
175 dataset pairs (**Supplementary Figure S3D and Figure S4**). Therefore, we
176 divided these pairs into two scenarios: zebrafish-excluded (139 pairs) and
177 zebrafish-associated (510 pairs) (**Methods**).

178 We compared the cell-typing performance of CAME with six benchmarking
179 methods including two marker-based methods SciBet [15] and Scamp [44], two
180 deep-learning methods Cell BLAST [17] and ItClust [18], one expression-based
181 method SingleCellNet [14], and one integration-based method Seurat-v3 [16]
182 in these two scenarios in terms of accuracy, macro-F1 score and weighted F1
183 score (**Methods**). Results showed that, in both scenarios, CAME distinctly
184 outperformed the others in most cases with statistical significance p-values
185 $< 10^{-16}$ and 10^{-54} using Wilcoxon signed-rank test for both zebrafish-
186 excluded and zebrafish-associated scenarios, respectively (**Figure 2A and B**,
187 **Supplementary Figures S5 and S6**).

188 To evaluate the robustness of CAME in the cases when the reference and
189 query datasets have inconsistent and insufficient sequencing depths, we
190 performed down-sampling experiments (at various sampling rates 75%, 50%,
191 25%, 10%) for read counts on the reference, query, and both reference and
192 query datasets. Again, CAME achieved superior performance compared to all
193 six benchmarking methods (**Figure 2C, Supplementary Figures S7 and S8**).
194 By contrast, when the down-sampling rates are extremely unbalanced, some
195 benchmarking methods may fail. For example, at a down-sampling rate of 0.1
196 for query datasets, Seurat detected too few anchors to abort integration for label
197 transfer and Scmap failed to find enough genes since the median expression
198 in the selected features is 0 in each cell cluster. All these results demonstrate
199 that CAME is robust to the insufficient and inconsistent sequencing depths
200 between reference and query pairs.

201 **CAME could robustly align homologous cell types across species and** 202 **multiple references**

203 In addition to the accurate cross-species cell-type assignment, CAME is also
204 capable of aligning homologous cell types from different species, even when
205 crossing distant species. For example, when aligning cell types between mouse
206 [29] and turtle [10], CAME successfully distinguished and aligned each major
207 type, like inhibitory and excitatory neurons, while the alignments by FastMNN,
208 Harmony, and Seurat were incapable. CAME also separated the neural
209 progenitor cells from excitatory neurons, while LIGER merged these two groups.
210 The visualization plots using Uniform Manifold Approximation and Projection
211 (UMAP) [31] of cell embeddings of Cell BLAST tend to lose some relations
212 between cell types, e.g., the inhibitory and excitatory neurons are not linearly

213 separable on the 2D plot (**Figure 3A**, and **Supplementary Figure S9**).

214 When handling multiple references and batch information is unavailable,
215 most integration methods will suffer from batch effects. In this situation, owing
216 to the semi-supervised manner, CAME can ignore the batch effects of reference
217 data. In contrast, other integration tools may suffer from diverse sources of
218 noises if the potential batch effects (such as noises from different individuals)
219 are not considered. For instance, when aligning human and mouse pancreas
220 cell types with human reference composed of eight batches, cells of the same
221 type but from different batches were still separated from each other. Besides,
222 the query cells tended to be “attracted” by reference cells of the same protocol
223 (**Figure 3B**). Even when the batch labels are given, for some of the
224 benchmarking methods (e.g., LIGER [19] and Seurat-v3 [16]), the reference
225 batch effects still existed after data integration (**Supplementary Figure S10**).

226 **CAME could accurately assign cell types in mouse brains and reveal cell-** 227 **type-specific gene modules**

228 We applied CAME to assign the major types of single cells from the primary
229 visual cortex and the anterior lateral motor cortex of mice [29], and used human
230 brain cells as the reference dataset [7], containing the cells from the hindbrain
231 that is not included in the mouse dataset. CAME achieved an accuracy of about
232 98%, so as Seurat and SciBet, superior to other benchmarking methods (94%
233 by ItClust, 93% by Cell BLAST, 92% by SingleCellNet, and only 55% by Scmap).
234 CAME also got a higher macro-F1 score (0.55) than that of Seurat (0.44) and
235 SciBet (0.46), indicating that CAME also accurately classified the small groups.
236 Specifically, those non-neuronal types accounting for a small proportion of
237 mouse cells were accurately assigned, including endothelial cells (accounting
238 for 0.6% of human cells and 0.85% of mouse cells) and its subclass, brain
239 pericytes (0.61% of human cells and 0.14% of mouse cells). The macrophages
240 (0.56% in mice) were classified as microglial cells (2.1% of human cells) that
241 are biologically similar to this type. Both oligodendrocyte precursor cells (OPC)
242 and oligodendrocytes in mice were originally assigned as oligodendrocytes
243 (0.75% of mouse cells) by the authors, but they were distinguished from each
244 other in the reference of the human data (**Figure 4A**). The identities of OPCs
245 were also verified by examining the expression of typical marker genes in each
246 cell type (**Figure 4B**). Besides, we found that the genes with top attentions from
247 each cell type showed high cell-type specificities, though these genes were
248 quite different across species (**Supplementary Figure S11A**).

249 Similar results were found when comparing four subtypes of the inhibitory
250 neurons (VIP+, SST+, LAMP5+, PVALB+) between humans and mice. CAME
251 still achieved a cell-typing accuracy of 98.3% and 95.5% for human-to-mouse
252 and mouse-to-human label transfers, respectively, which are consistently
253 higher than that of the benchmarking methods (93.4% and 92.0% for SciBet,
254 84.3% and 51.3% for SingleCellNet, 98.0% and 78.9% for Cell BLAST, 97.3%
255 and 87.3% for ItClust, 69.5% and 78.9% for Scmap, 94.2% and 87.2% for

256 Seurat) (**Supplementary Figure S12A and B**), although differentially
257 expressed genes (DEGs) for each homologous subtype seems not
258 transferrable across species (**Supplementary Figure S12C and D**). The
259 UMAP plots of cell embeddings showed that these major homologous cell-types
260 were well aligned with each other. This suggested that the major types of brain
261 cells in humans and mice are well conserved (**Supplementary Figure S9**).

262 CAME also gave interpretable gene embeddings and enabled us to explore
263 both intra- and inter-species relationships between genes. The UMAP plots of
264 gene embeddings showed that the relative positions of human and mouse
265 homologous genes were very consistent (**Figure 5C**). We further demonstrated
266 the averaged gene expression profile on the UMAP plots of gene embeddings,
267 where each point represents a gene (**Figure 4C** and **Supplementary Figure**
268 **S11B**). It is worth noting that the neighbor genes tend to be co-expressed in the
269 same cell types, such as those in excitatory, inhibitory neurons,
270 oligodendrocytes, and OPCs (**Figure 4D**). There were more cell type-specific
271 genes in human oligodendrocytes than in mice, indicating the evolutionary
272 divergence between humans and mice. A population of genes was only
273 detected in the human dataset, and most of them were associated with Purkinje
274 cells and cerebellum granule cells, which were not detected in the mouse
275 dataset due to their sources from different brain regions. These genes were
276 arranged where there were few mouse genes around (**Figure 4F**, and
277 **Supplementary Figure S11B**).

278 The aligned gene embeddings across species can facilitate us to jointly
279 extract cell type-specific gene modules with different degrees of conservancies
280 between species, and each module corresponds to a cell type like OPCs, or
281 related cell types like endothelial cells and its subtypes (**Figure 4E** and
282 **Methods**). As expected, based on gene ontology (GO) [33, 34] enrichment
283 analysis, we found that the functions associated with most homologous gene
284 modules were generally consistent with each other (**Supplementary Table S2**).
285 For example, both the human and mouse genes in module 2 (which was
286 associated with inhibitory neurons) tended to relate functions like “forebrain
287 neuron differentiation” and “learning or memory”. Both the human and mouse
288 genes in module 6 (corresponding human microglia and mouse macrophage)
289 were related to functions like “positive regulation of cytokine production”, and
290 “leukocyte migration”. By contrast, the function “ventral spinal cord
291 development” was only enriched in human module 3 but not in mice,
292 considering their gene members were quite different; though they were both
293 associated with the function “cell differentiation in hindbrain” and “cerebellar
294 cortex formation”.

295 **CAME could reveal conservative expression dynamics during** 296 **spermatogenesis between human and macaque**

297 Comparison of continuous biological processes between two species is of much
298 interest in evolutionary biology. We applied CAME to two scRNA-seq datasets

299 from human and macaque testicular single cells [9] with the former as the
300 reference one. CAME achieved a very distinct cell-typing accuracy of 95.0%
301 (86.0% for SciBet, 89.2% for SingleCellNet, 76.1% for Cell BLAST, 53.4% for
302 ItClust, 87.3% for Scmap, 89.1% for Seurat), and a precise alignment of the
303 homologous cell types of human and macaque with each other (**Figure 5A and**
304 **B**). Besides, the labeled spermatogonia, spermatocyte, round spermatid, and
305 elongating cells are correctly merged along the underlying differentiation
306 trajectory. This suggested that CAME could well decipher the conserved four-
307 stage spermatogenesis processes of humans and macaques.

308 Very interestingly, the continuously dynamic changing process of
309 spermatogenesis can also be revealed by the UMAP plot of gene embeddings
310 (**Figure 5C**). As illustrated, CAME extracted four sets of genes, including some
311 typical marker ones [32], that are highly co-expressed in the four main stages
312 of spermatogenesis and form well-organized expression dynamics, suggesting
313 the order of critical gene activations during spermatogenesis (**Figure 5C**). By
314 joint extraction of gene modules, we found that the four stages of
315 spermatogenesis were quite conservative from the aspect of gene modules
316 (**Figure 5D and E**). For example, modules 3, 4, and 0 were highly expressed
317 in spermatogonia and spermatocyte respectively for both humans and
318 macaques. And round spermatids and elongating spermatids shared modules
319 2, 1, and 5 in different degrees. Typically, both human and macaque module 4
320 was associated with functions like “RNA splicing”, and module 1 was associated
321 with “sperm motility” and “spermatid development/differentiation”, which were
322 typical characteristics of elongating spermatids (**Supplementary Table S3**).

323 **Conclusions**

324 Cross-species comparative and integrative analysis at single-cell resolution has
325 deepened our understanding of the origin and evolutionary mechanisms of
326 cellular states. Exploring the conservative and divergent characteristics of
327 homologous cell states between human and other model and non-model
328 species, for example, can help us to determine the animal model for studying
329 human disease [5-7].

330 However, in addition to technical noises, the systematic shift of gene
331 expressions associated with distinct species and the uncertainty of the
332 orthologous genes make it much more difficult than within-species data
333 integration. Moreover, existing approaches for cross-species integration were
334 mainly based on one-to-one homologous genes. However, when it is needed
335 to align cell types across long distant species, especially when a large number
336 of gene duplications were involved during the evolution process [27,28],
337 considering only the one-to-one homologous genes will inevitably lose a lot of
338 important information. Even so, cells of homologous types are thought to have
339 similar expression patterns, that is, they may co-express a cell type-specific
340 combination of genes. These genes may not be easy to be identified as the

341 marker genes with high expression levels but can act as “bridges” between cells
342 that co-expressed them. Besides, the gene-homology mappings can bridge the
343 gene nodes of two species, where the non-one-to-one homologies can also be
344 used.

345 Thus, we take the gene expression matrix as a bipartite graph with cell and
346 gene nodes and utilize the gene homologies to form a multipartite graph. Based
347 on this, we proposed CAME to utilize a cell-gene heterogeneous graph neural
348 network to facilitate the “message-passing” from one species to the other.
349 CAME can achieve the alignment of both cells and genes from different species.
350 As a result, CAME can not only achieve accurate and robust cell-type
351 assignment, but also reveal biological insights into the conservative and
352 divergent characteristics between species. When handling multiple references,
353 most integration approaches have to perform pairwise alignment for individual
354 batches, where the order of pairwise alignment can affect the results and the
355 computational complexity rises quadratically with the number of batches.
356 Others like Harmony [30] and Cell BLAST [17] are capable to align multiple
357 datasets simultaneously. We demonstrated that CAME can remove batch
358 effects for multiple references even when batch labels are not provided. This is
359 an important characteristic for integrating various datasets and constructing a
360 unified cell-typing reference.

361 It should be noticed that the heterogeneous graph neural network structure
362 of CAME can also be applied to the scenario of within-species data integration,
363 or when we consider only the one-to-one homologous genes. The only
364 adjustment is to replace each gene-gene edge with a single gene node.
365 Moreover, this strategy can be applied for multi-omics label transfer and data
366 integration. In summary, we believe that CAME will serve as a powerful tool for
367 integrative and comparative analysis across species as well as multi-omics
368 integration.

369

370 Methods

371 Build a heterogeneous cell-gene graph

372 Let's denote a gene expression matrix with N cells and M genes as $X =$
 373 $(X_1, X_2, \dots, X_N)^T \in \mathbb{R}^{N \times M}$, where each row $X_i = (x_{i1}, x_{i2}, \dots, x_{iM}) \in \mathbb{R}^M$ with an
 374 element x_{ij} representing the (normalized) expression value of a cell i in a
 375 gene j . We take $X^{(R)} \in \mathbb{R}^{N_R \times M_R}$ and $X^{(Q)} \in \mathbb{R}^{N_Q \times M_Q}$ as the reference and
 376 query datasets respectively, $Y = (y_1, y_2, \dots, y_{N_R}) \in \mathbb{R}^{N_R}$ as the cell-type labels
 377 of the reference dataset and a set of gene pairs $\{(g_i, g_j)\}_{ij}$ to indicate the
 378 homology between two species. Note that M_R is not necessarily equal to M_Q .

379 The reference and query expression matrices and the homology together are
 380 represented as a heterogeneous cell-gene graph with each node acting as a
 381 cell or a gene (**Figure 1A**). A cell-gene edge in the graph indicates that this cell
 382 has non-zero expression of the gene, a gene-gene edge indicates a homology
 383 between each other, and a cell-cell edge indicates the expression profiles of
 384 these two cells are similar to each other. In other words, in this graph, there are
 385 two types of nodes, cell and gene, and six types of edges (relations) including
 386 "a cell expresses a gene", "a gene is expressed by a cell", "cell-cell similarity",
 387 "gene-gene homology", "cell self-loop" and "gene self-loop".

388 Design a heterogeneous graph neural network

389 CAME adopts a heterogeneous graph neural network, which was motivated by
 390 a relational graph convolutional network [24] for a graph of homogeneous
 391 nodes but heterogeneous edges. We denote the convolution weights for these
 392 six edge types as W_{cg} , W_{gc} , W_{cc} , W_{gg} , W_c and W_g , respectively (**Figure 1B**).
 393 For each cell i , its initial embedding (the 0-th layer) is calculated as:

$$394 \quad h_{c_i}^{(0)} = \sigma(W_c^{(0)} x_{c_i} + b_c^{(0)}),$$

395 where σ is the ReLU activation function; x_{c_i} is the gene expressions in the
 396 cell i (one-to-one homologous genes are taken as the common input features)
 397 and $b_c^{(0)} \in \mathbb{R}^{d^{(0)}}$ is the learnable bias vector. The genes, however, lack the
 398 initial embeddings in the 0-th layer and can be aggregated from their neighbor
 399 cells as follows:

$$400 \quad h_{g_j}^{(0)} = \sigma\left(\sum_{i \in \mathcal{N}_{g_j}^c} \frac{1}{z_{g_j, c}} W_{cg}^{(0)} x_{c_i} + b_g^{(0)}\right),$$

401 where $\mathcal{N}_{g_j}^c$ is the set of cells that have expressed the gene j , and $z_{g_j, c} = |\mathcal{N}_{g_j}^c|$

402 is the normalization factor. This approach keeps the number of model
 403 parameters stay constant to the number of genes, which differs from the
 404 commonly used initialization that assigns a learnable embedding for those
 405 nodes without input features, where the increasing number of model
 406 parameters might lead to an overfitted model. It can also allow inductive
 407 learning for the genes not involved in the training process.

408 While in each hidden layer $l \geq 1$, the node features for the cell i and the
 409 gene j can be calculated as:

$$410 \quad h_{c_i}^{(l)} = \sigma \left(\sum_{j \in \mathcal{N}_{c_i}^g} \frac{1}{z_{c_i, g}} W_{gc}^{(l)} h_{g_j}^{(l-1)} + \sum_{k \in \mathcal{N}_{c_i}^c} \frac{1}{z_{c_i, c}} W_{cc}^{(l)} h_{c_k}^{(l-1)} + W_c^{(l)} h_{c_i}^{(l-1)} + b_c^{(l)} \right),$$

411 and

$$412 \quad h_{g_j}^{(l)} = \sigma \left(\sum_{i \in \mathcal{N}_{g_j}^c} \frac{1}{z_{g_j, c}} W_{cg}^{(l)} h_{c_i}^{(l-1)} + \sum_{k \in \mathcal{N}_{g_j}^g} \frac{1}{z_{g_j, g}} W_{gg}^{(l)} h_{g_k}^{(l-1)} + W_g^{(l)} h_{g_j}^{(l-1)} + b_g^{(l)} \right),$$

413 respectively. Note that we treat the edges between homologous genes and the
 414 self-loop on each gene identically, i.e., $W_{gg}^{(l)} = W_g^{(l)}$. To boost the ‘message’
 415 flow between reference and query nodes, we adopt a recurrent convolution,
 416 where the parameters are shared across the hidden layers, that is, $W_{gc}^{(l)} =$
 417 $W_{gc}, W_{cg}^{(l)} = W_{cg}, W_{gg}^{(l)} = W_g^{(l)} = W_g, W_c^{(l)} = W_c$ and $b_c^{(l)} = b_c, b_g^{(l)} = b_g$ for $1 \leq$
 418 $l \leq L$, where L is the total number of the hidden layers. We recommend to set
 419 L as 2 or 3 in practice, and the default setting is 2. We also adopt the layer
 420 normalization for all the hidden states to facilitate fast training convergence and
 421 high performance (**Supplementary Figure S1**).

422 When it comes to the cell-type classifier, we adopt the attention mechanism
 423 for graph convolution [25], where each cell pays distinct attention to its neighbor
 424 genes. Specifically, for each cell i , the output states $h_{c_i}^{out}$ for cell-type
 425 identification is aggregated from their neighbor genes:

$$426 \quad h_{c_i}^{out} = \sum_{j \in \mathcal{N}_{c_i}^g} \alpha_{ij} W_g^{out} h_{g_j}^{(L)} + b^{out},$$

427 where α_{ij} is the attention that the cell i pays to the gene j , calculated as:

$$428 \quad \alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}_{c_i}^g} \exp(e_{ik})},$$

429 with

$$430 \quad e_{ij} = \text{leakyReLU} \left(a^T \left[W_c^{out} h_{c_i}^{(L)} \parallel W_g^{out} h_{g_j}^{(L)} \right] \right).$$

431 In addition, we use multi-head attention to enhance the model capacity and
 432 robustness, where there are several attention-heads with their own parameters,
 433 and their outputs are merged by taking averages:

$$434 \quad h_{c_i}^{out} = \frac{1}{K} \sum_{k=1}^K h_{c_i}^{out,k} = \frac{1}{K} \sum_{k=1}^K \left(\sum_{j \in \mathcal{N}_{c_i}^g} \alpha_{ij} W_g^{out,k} h_{g_j}^{(L)} + b^{out,k} \right),$$

435 where K is the total number of attention-heads, set as 8 by default.

436 Finally, the output layer states for cell-type classification were normalized in
 437 two different ways: (1) the *softmax* function over cell types for multi-class
 438 classification:

$$439 \quad Y' = \text{softmax}(H^{out}), \quad H^{out} = (h_{c_1}^{out}, \dots, h_{c_N}^{out})^T,$$

440 where $Y' \in \mathbb{R}^{N \times T}$ and each row is the predicted probabilities over the T cell
 441 types for a cell; (2) the sigmoid function for multi-label classification:

$$442 \quad Y'' = \text{sigmoid}(H^{out}) = \frac{1}{1 + \exp(H^{out})},$$

443 where $Y'' \in \mathbb{R}^{N \times T}$ and each element Y''_{it} is the predicted probability of the cell
 444 type t for the cell i .

445 **The classification loss and label smoothing**

446 The classification loss for cells in reference datasets is calculated by the
 447 weighted cross-entropy loss combined with L_2 regularization as below:

$$448 \quad L_c(\mathbf{X}_R, \mathbf{Y}_R) = \frac{1}{N_R} \sum_{i=1}^{N_R} \left[\sum_{t=1}^T w_t Y_{it} \ln(Y'_{it}) + \sum_{t=1}^T w_t Y_{it} \ln(Y''_{it}) \right] + \lambda \|\theta\|_2^2$$

$$449 \quad = \frac{1}{N_R} \sum_{i=1}^{N_R} \sum_{t=1}^T w_t Y_{it} \ln(Y'_{it} Y''_{it}) + \lambda \|\theta\|_2^2,$$

450 where w_t is the class-weight for cell-type t satisfying $\sum_{t=1}^T w_t = 1$. To avoid
 451 the model being dominated by the major populations and ignoring those rare
 452 types, we set $w_t \propto \frac{1}{\sqrt{N_t}}$ and N_t is the number of cells of cell type t in the
 453 reference dataset. θ represents all the learnable parameters and λ_1 is the
 454 penalization coefficient that controls the power of L_2 regularization, and the
 455 default value of λ_1 is 0.01.

456 To prevent the model from being overconfident and improve the stability and
 457 generalization of the model, we utilize label smoothing [35]. We minimize the
 458 cross-entropy between the modified targets $Y^{LS} \in \mathbb{R}^{N_R \times T}$ and the model
 459 outputs Y' , where $Y_{it}^{LS} = Y_{it}(1 - \alpha) + \alpha/K$, and the final objective function is
 460 as below:

461
$$L_{sc} = (1 - \epsilon)L_c + \frac{\epsilon}{T} \sum_{t=1}^T \frac{1}{N_R} \sum_{i=1}^{N_R} \ln(Y'_{it} Y''_{it}),$$

462 where ϵ controls the degree of smoothness, set as 0.1 by default. Finally,
463 CAME adopts Adam optimizer [36] with a learning rate of 0.001 for training.

464 **Checkpoint selection**

465 When training the heterogeneous graph neural network, we would like to
466 choose the epoch where the classification result of query datasets achieves the
467 highest accuracy. However, in practice, the exact type labels of the query cells
468 are unknown, hindering us from choosing the best model. We put forward a
469 metric to approximate the accuracy. Specifically, we first cluster the query cells
470 to get the pseudo-labels $Y^{cluster}$ for the query cells and introduce adjusted
471 mutual information (AMI) [37] to account for the chance between the model-
472 predicted cell-type labels and the pseudo-labels of the query cells to help
473 decide when to stop. AMI is defined as

474
$$AMI(Y^{cluster}, Y') = \frac{MI(Y^{cluster}, Y') - E[MI(Y^{cluster}, Y')]}{\text{mean}\{H(Y^{cluster}), H(Y')\} - E[MI(Y^{cluster}, Y')]},$$

475 where $H(X)$ is the entropy of X , $MI(X, Y)$ is the mutual information between
476 variables X and Y . $E[MI(Y^{cluster}, Y')]$ is the expected mutual information
477 based on a “permutation model” [38], in which cluster labels are generated
478 randomly subject to having a fixed number of clusters and points in each cluster.
479 We think that a well-trained model is expected to preserve the intrinsic data
480 structure so that the predicted labels should be highly consistent with the
481 pseudo-labels to some extent. We run the model with 400 epochs and choose
482 the checkpoint with the largest AMI. The clustering process will be described in
483 the section “*pre-clustering of the query cells*” in detail.

484 **Training using the mini-batches on sub-graphs**

485 When training CAME on the graphics processing unit (GPU), the size of a
486 dataset will be limited by the GPU memory. For example, training CAME on
487 100,000 cells could take about 13.75GB of memory, which exceeds the graphic
488 memory of most GPUs. To handle this issue, we utilized a mini-batch training
489 process by using the graph segmentation technique. Specifically, we first
490 randomly divided all the cells (including cells in reference and query) into
491 several groups, taken as mini-batches. For each mini-batch, we created a node-
492 induced subgraph for a given group of cells, which contains all the cells in this
493 group and all the genes expressed by these cells. Then, we iterated all
494 subgraphs and feed the subgraphs to the graph neural network one by one. All
495 the parameters will be updated for each mini-batch training process. We
496 performed extensive experiments by using mini-batch training process and
497 found it is suitable to choose batch-size as 8192 or more, for that it achieved
498 the comparable accuracy compared with the whole graph training

499 **(Supplementary Figure S14A)** and the cost of GPU memory stays constant
500 (2.4GB) for datasets at different scales **(Supplementary Figure S14B)**. Such
501 low consumption of graphic memory means you can use CAME on almost all
502 graphics cards. It is worth noting that the runtime of the batch-training process
503 will be largely increased **(Supplementary Figure S14B)** since we cannot feed
504 forward the whole graph on a single epoch.

505 **Preprocessing of the single-cell datasets**

506 For each scRNA-seq dataset, we first normalized the counts of each cell by its
507 library size with a scale factor multiplied (10,000 by default) and log-
508 transformed with a pseudo-count added for the downstream analysis.

509 **Gene selection**

510 Highly variable genes (HVGs) and differentially expressed genes (DEGs) are
511 generally thought to be highly informative and the latter is especially useful for
512 cell-type characterization. Therefore, we used both HVGs and DEGs and
513 extended them using homologous mappings to form the highly informative gene
514 (HIG) sets for constructing the heterogeneous graph. We adopted the same
515 approach as used in Seurat-v2 [39] with ScanPy [40] built-in function
516 `highly_variable_genes()` to identify HVGs, separately from both reference and
517 query data. Specifically speaking, it calculated the average expression and
518 dispersion (variance/mean) for each gene and placed these genes into several
519 bins based on the (log-transformed) average expression. The normalized
520 dispersions were then obtained by scaling with the mean and standard
521 deviation of the dispersions within each bin. We selected the top 2000 genes
522 with the highest dispersions as HVGs of that dataset. We computed the DEGs
523 separately for reference and query dataset by Student's t-test, which is done
524 through `rank_genes_groups()` function from the ScanPy package [40]. For
525 reference data, cells are grouped by their cell-type labels, while for the query
526 data, cells are grouped by their pseudo-labels, i.e., the pre-clustering labels.

527 Genes used as the cell-node features should be shared between species (or
528 datasets). For both reference and query datasets, we first took the top 50 DEGs
529 for each cell group and retained genes with one-to-one homology in the other
530 species. We then took the union of the resulting two sets of genes for input. The
531 resulting number of genes used for defining cell-node features ranges from 240
532 to 400 for distant species pairs (human to zebrafish for example) and from 400
533 to 900 for the others.

534 We combined both HVGs and DEGs from reference and query data to decide
535 the node genes used for training the graph neural network. Specifically, we first
536 took the union of the HVGs and DEGs for each dataset, denoted as \mathcal{G}_r and \mathcal{G}_q
537 for reference and query respectively. Then we extracted the genes that have
538 homologies in \mathcal{G}_r from the query data, and the homologous genes for \mathcal{G}_q from

539 the reference data denoted as $\mathcal{G}_r^{(homo)}$ and $\mathcal{G}_q^{(homo)}$ respectively. Finally, we

540 determined $\mathcal{G}_r \cup \mathcal{G}_q^{(homo)}$, the union of \mathcal{G}_r and $\mathcal{G}_q^{(homo)}$, as the node genes for
541 the reference species and $\mathcal{G}_r^{(homo)} \cup \mathcal{G}_q$ as the node genes for the query
542 species. The tables containing gene homology information for each species pair
543 were downloaded from the BioMart web server
544 (<http://www.ensembl.org/biomart/martview>) [41].

545 **Construction of the single-cell graphs based on KNNs**

546 The normalized expression matrices were centralized and scaled within each
547 dataset, followed by principal component analysis (PCA) to reduce the
548 dimensionality. We searched approximate KNNs for each cell based on the top
549 30 PCs with the highest explained variances. We adopted $k = 5$ neighbors for
550 each cell to make the graph sparse enough for computational efficiency. These
551 neighbor connections provided “cell-cell” edges as a part of the heterogeneous
552 graph.

553 **Pre-clustering of the query cells**

554 To facilitate model selection, we pre-clustered the query cells using a graph-
555 based clustering method, that is, performing community detection using the
556 Leiden algorithm [42] on the single-cell KNN graph. We constructed the KNN
557 graph in almost the same way as described above, except that the number of
558 neighbors k was set as 20 and the clustering resolution is set as 0.4 by default.

559 **Unifying cell-type labels across datasets**

560 For data downloaded from the Cell BLAST web server [17], the cell-type labels
561 were already unified by Cell Ontology [43], a structured vocabulary for cell types.
562 While for unifying annotations from the other datasets, we directly referred to
563 Cell Ontology and manually adjusted the annotations. The annotations were
564 used as ground truth.

565 **Gene module extraction**

566 To extract cell type-specific gene modules shared between species, we took all
567 the gene embeddings (of both species) on the last hidden layer and performed
568 KNN searching for each gene. Like clustering cells, we performed Leiden
569 community detection on the KNN graph of genes. The clustering resolution was
570 set as 0.8 by default.

571 **Calculating weights between gene modules**

572 The weights S_{ij} between homologous gene modules Mod_i and Mod_j on the
573 abstracted graph were calculated as follows:

$$574 \quad S_{ij} = \frac{\sum_{(g_1 \in Mod_i) \wedge (g_2 \in Mod_j)} sim(h_{g_1}, h_{g_2})}{\max(|Mod_i|, |Mod_j|)},$$

575 where h_g is the embedding vector of gene g and $sim(\cdot, \cdot)$ is the similarity
576 function, cosine similarity by default. $|Mod|$ represents the number of genes in
577 this module.

578 **Benchmarking cell-type assignment**

579 For benchmarking cell-type assignment, we collected 54 scRNA-seq datasets
580 from five tissues across seven different species (**Supplementary Figure 3A,**
581 **Supplementary Table 1**), paired datasets of different species within the same
582 tissue, and filtered those pairs where more than 50% of query cells are
583 unresolved in the reference cell types, resulting 649 cross-species dataset pairs.
584 For each dataset, we removed the cell types of less than 10 cells. CAME was
585 compared with six benchmarking methods including Seurat V3 [16], ItClust [18],
586 Scmap [44], SingleCellNet [14], SciBet [15], and Cell BLAST [17]. For Seurat
587 V3, we input the raw data, used the default normalize process by
588 `NormalizeData()` function, extracted the top 2000 HVGs by its
589 `FindVariableFeatures()` function for reference and query respectively, and
590 performed further annotation process as described in its documentation. For
591 ItClust, since it provides an automatic workflow including preprocessing and
592 annotation, we input the raw data. For Scmap, we log-transformed the raw
593 counts with pseudo-count 1 added and used its inherited function
594 `selectFeaures()` to select the top 2000 HVGs with a threshold=0.1 in function
595 `scmapCluster()` (which works better for the cross-species scenario than its
596 default value). For SingleCellNet, we also input raw data as it suggested, used
597 `splitCommon` function to split for training and assessment, employed `expTMraw`
598 function to transform training data, and then used `scn_predict` to make
599 predictions for query dataset. For SciBet, we used R to perform all the
600 operations. We first input the library-size-normalized data calculated by `cpm()`
601 function of package `edgeR` [45] and used `SelectGene_R()` function from `SciBet`
602 package to select 2000 HVGs, and used `SciBet_R()` function to annotate the
603 query data. For Cell BLAST, we used the raw data as input and used
604 `find_variable_genes()` to select HVGs with default parameters and took the
605 union of the HVGs between reference and query datasets. After that, the
606 datasets were combined together to remove their batch effects by using
607 function `fit_DIRECTi()` with `lambad_reg=0.001` as suggested by the original
608 authors to stabilize the training process. Cell BLAST also provides a supervised
609 training process that leverages the cell type labels of reference datasets to
610 perform label transfer. However, it led to a 4% decrease in the average
611 accuracy compared with their previous batch effects correction process. All
612 hyper-parameters not mentioned were set with default values in these six
613 packages.

614 To evaluate the performance of the cell-type assignment, we adopted three
615 metrics: Accuracy, MarcoF1, and WeightedF1. Accuracy is the most common
616 criterion and it directly measures how many of the predictions are the same as
617 the actual ones:

618
$$Acc = \frac{\#\{Y'==Y\}}{\#\{Y\}},$$

619 where # is the sign of cardinality. Specifically, $\#\{Y_{true}\}$ means the number of
620 the total cells and $\#\{Y' == Y\}$ means the number of correctly predicted ones.

621 We also used *MacroF1* and *WeightedF1* which consider the F_1 -score for
622 each cell type. For a binary classification task, precision and recall are
623 calculated as

624
$$precision = \frac{TP}{TP+FP},$$

625 and

626
$$recall = \frac{TP}{TP+FN},$$

627 respectively, where TP, FP, and FN represent the number of true positives, false
628 positives, and false negatives, respectively.

629 The F_1 -score is the harmonic mean of *precision* and *recall*:

630
$$F_1 = \frac{2 \times precision \times recall}{precision + recall},$$

631 and the *MacroF1* is defined as the average of class-wise F_1 -scores,

632
$$MacroF_1 = \frac{1}{T} \sum_{c=1}^t F_1^{(t)},$$

633 where $F_1^{(t)}$ represents the F_1 -score for cell type t . The *WeightedF1*
634 considers the proportion of each class,

635
$$WeightedF_1 = \sum_{t=1}^T \frac{N_t}{N} \times F_1^{(t)},$$

636 where N_t/N represents the proportion of type t in all cells.

637 **Benchmarking data integration**

638 FastMNN [46], Harmony [30], and Seurat-v3 [16] were performed using the
639 corresponding R package through SeuratWrapper, following the online
640 documents with default settings. FastMNN, Harmony, and Seurat shared the
641 same normalization and the top 2000 HVGs by Seurat function NormalizeData
642 and FindVariableFeatures, respectively. Harmony was performed on the PCA-
643 reduced embeddings. The number of reduced dimensions for these three
644 methods was set as 50 for all pairs of datasets. LIGER [19] took the raw count
645 data as input and run with the default pipeline. Cell BLAST [17] was performed
646 using its Python package, following the standard pipeline with the default
647 settings.

648

649 Acknowledgments

650 This work has been supported by the National Key Research and Development
651 Program of China [2019YFA0709501]; the Strategic Priority Research Program of the
652 Chinese Academy of Sciences (CAS) [XDPB17], the Key-Area Research and
653 Development of Guangdong Province [2020B1111190001], the National Natural
654 Science Foundation of China [61621003]; the National Ten Thousand Talent Program
655 for Young Top-notch Talents, the CAS Frontier Science Research Key Project for Top
656 Young Scientist [QYZDB-SSW-SYS008], and the Shanghai Municipal Science and
657 Technology Major Project [2017SHZDZX01].
658

659 References

- 660 [1] E.Z. Macosko, A. Basu, R. Satija, J. Nemesh, K. Shekhar, M. Goldman, I. Tirosh,
661 A.R. Bialas, N. Kamitaki, E.M. Martersteck, J.J. Trombetta, D.A. Weitz, J.R. Sanes,
662 A.K. Shalek, A. Regev, S.A. McCarroll, Highly Parallel Genome-wide Expression
663 Profiling of Individual Cells Using Nanoliter Droplets, *Cell* 161(5) (2015) 1202-1214.
664 [2] A.A. Kolodziejczyk, J.K. Kim, V. Svensson, J.C. Marioni, S.A. Teichmann, The
665 technology and biology of single-cell RNA sequencing, *Molecular cell* 58(4) (2015)
666 610-20.
667 [3] J.C. Marioni, D. Arendt, How Single-Cell Genomics Is Changing Evolutionary and
668 Developmental Biology, *Annual Review of Cell and Developmental Biology* 33(1)
669 (2017) 537-553.
670 [4] M.E.R. Shafer, Cross-Species Analysis of Single-Cell Transcriptomic Data, *Front*
671 *Cell Dev Biol* 7 (2019) 175-175.
672 [5] E. Drokhlyansky, C.S. Smillie, N. Van Wittenberghe, M. Ericsson, G.K. Griffin, G.
673 Eraslan, D. Dionne, M.S. Cuoco, M.N. Goder-Reiser, T. Sharova, O. Kuksenko, A.J.
674 Aguirre, G.M. Boland, D. Graham, O. Rozenblatt-Rosen, R.J. Xavier, A. Regev, The
675 Human and Mouse Enteric Nervous System at Single-Cell Resolution, *Cell* 182(6)
676 (2020) 1606-1622.e23.
677 [6] L. Geirsdottir, E. David, H. Keren-Shaul, A. Weiner, S.C. Bohlen, J. Neuber, A.
678 Balic, A. Giladi, F. Sheban, C.-A. Dutertre, C. Pfeifle, F. Peri, A. Raffo-Romero, J.
679 Vizioli, K. Matiasek, C. Scheiwe, S. Meckel, K. Mätz-Rensing, F. van der Meer, F.R.
680 Thormodsson, C. Stadelmann, N. Zilkha, T. Kimchi, F. Ginhoux, I. Ulitsky, D. Erny, I.
681 Amit, M. Prinz, Cross-Species Single-Cell Analysis Reveals Divergence of the
682 Primate Microglia Program, *Cell* 179(7) (2019) 1609-1622.e16.
683 [7] R.D. Hodge, T.E. Bakken, J.A. Miller, K.A. Smith, E.R. Barkan, L.T. Graybuck, J.L.
684 Close, B. Long, N. Johansen, O. Penn, Z. Yao, J. Eggermont, T. Höllt, B.P. Levi, S.I.
685 Shehata, B. Aevermann, A. Beller, D. Bertagnolli, K. Brouner, T. Casper, C. Cobbs, R.
686 Dalley, N. Dee, S.L. Ding, R.G. Ellenbogen, O. Fong, E. Garren, J. Goldy, R.P.
687 Gwinn, D. Hirschstein, C.D. Keene, M. Keshk, A.L. Ko, K. Lathia, A. Mahfouz, Z.
688 Maltzer, M. McGraw, T.N. Nguyen, J. Nyhus, J.G. Ojemann, A. Oldre, S. Parry, S.
689 Reynolds, C. Rimorin, N.V. Shapovalova, S. Somasundaram, A. Szafer, E.R.
690 Thomsen, M. Tieu, G. Quon, R.H. Scheuermann, R. Yuste, S.M. Sunkin, B.

- 691 Lelieveldt, D. Feng, L. Ng, A. Bernard, M. Hawrylycz, J.W. Phillips, B. Tasic, H. Zeng,
692 A.R. Jones, C. Koch, E.S. Lein, Conserved cell types with divergent features in
693 human versus mouse cortex, *Nature* 573(7772) (2019) 61-68.
- 694 [8] A. Seb e-Pedr s, E. Chomsky, K. Pang, D. Lara-Astiaso, F. Gaiti, Z. Mukamel, I.
695 Amit, A. Hejnol, B.M. Degnan, A. Tanay, Early metazoan cell type diversity and the
696 evolution of multicellular gene regulation, *Nat Ecol Evol* 2(7) (2018) 1176-1188.
- 697 [9] A.N. Shami, X. Zheng, S.K. Munyoki, Q. Ma, G.L. Manske, C.D. Green, M.
698 Sukhwani, K.E. Orwig, J.Z. Li, S.S. Hammoud, Single-Cell RNA Sequencing of
699 Human, Macaque, and Mouse Testes Uncovers Conserved and Divergent Features
700 of Mammalian Spermatogenesis, *Developmental Cell* (2020).
- 701 [10] M.A. Tosches, T.M. Yamawaki, R.K. Naumann, A.A. Jacobi, G. Tushev, G.
702 Laurent, Evolution of pallium, hippocampus, and cortical cell types revealed by
703 single-cell transcriptomics in reptiles, *Science (New York, N.Y.)* 360(6391) (2018)
704 881-888.
- 705 [11] J. Wang, H. Sun, M. Jiang, J. Li, P. Zhang, H. Chen, Y. Mei, L. Fei, S. Lai, X.
706 Han, X. Song, S. Xu, M. Chen, H. Ouyang, D. Zhang, G.-C. Yuan, G. Guo, Tracing
707 cell-type evolution by cross-species comparison of cell atlases, *Cell Reports* 34(9)
708 (2021) 108803.
- 709 [12] A.W. Zhang, C. O'Flanagan, E.A. Chavez, J.L.P. Lim, N. Ceglia, A. McPherson,
710 M. Wiens, P. Walters, T. Chan, B. Hewitson, D. Lai, A. Mottok, C. Sarkozy, L. Chong,
711 T. Aoki, X. Wang, A.P. Weng, J.N. McAlpine, S. Aparicio, C. Steidl, K.R. Campbell,
712 S.P. Shah, Probabilistic cell-type assignment of single-cell RNA-seq for tumor
713 microenvironment profiling, *Nature Methods* 16(10) (2019) 1007-1015.
- 714 [13] X. Shao, J. Liao, X. Lu, R. Xue, N. Ai, X. Fan, scCATCH: Automatic Annotation
715 on Cell Types of Clusters from Single-Cell RNA Sequencing Data, *iScience* 23(3)
716 (2020) 100882.
- 717 [14] Y. Tan, P. Cahan, SingleCellNet: A Computational Tool to Classify Single Cell
718 RNA-Seq Data Across Platforms and Across Species, *Cell Systems* 9(2) (2019) 207-
719 213.e2.
- 720 [15] C. Li, B. Liu, B. Kang, Z. Liu, Y. Liu, C. Chen, X. Ren, Z. Zhang, SciBet as a
721 portable and fast single cell type identifier, *Nature Communications* 11(1) (2020)
722 1818.
- 723 [16] T. Stuart, A. Butler, P. Hoffman, C. Hafemeister, E. Papalexi, W.M. Mauck, 3rd, Y.
724 Hao, M. Stoeckius, P. Smibert, R. Satija, Comprehensive Integration of Single-Cell
725 Data, *Cell* 177(7) (2019) 1888-1902.e21.
- 726 [17] Z.-J. Cao, L. Wei, S. Lu, D.-C. Yang, G. Gao, Searching large-scale scRNA-seq
727 databases via unbiased cell embedding with Cell BLAST, *Nature Communications*
728 11(1) (2020) 3458.
- 729 [18] J. Hu, X. Li, G. Hu, Y. Lyu, K. Susztak, M. Li, Iterative transfer learning with
730 neural network for clustering and cell type classification in single-cell RNA-seq
731 analysis, *Nature Machine Intelligence* 2(10) (2020) 607-618.
- 732 [19] J.D. Welch, V. Kozareva, A. Ferreira, C. Vanderburg, C. Martin, E.Z. Macosko,
733 Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell
734 Identity, *Cell* 177(7) (2019) 1873-1887.e17.

- 735 [20] L. Zhang, S. Zhang, Learning common and specific patterns from data of
736 multiple interrelated biological scenarios with matrix factorization, *Nucleic acids*
737 *research* 47(13) (2019) 6606-6617.
- 738 [21] D. Arendt, J.M. Musser, C.V.H. Baker, A. Bergman, C. Cepko, D.H. Erwin, M.
739 Pavlicev, G. Schlosser, S. Widder, M.D. Laubichler, G.P. Wagner, The origin and
740 evolution of cell types, *Nature reviews. Genetics* 17(12) (2016) 744-757.
- 741 [22] S. Aibar, C.B. González-Blas, T. Moerman, V.A. Huynh-Thu, H. Imrichova, G.
742 Hulselmans, F. Rambow, J.C. Marine, P. Geurts, J. Aerts, J. van den Oord, Z.K. Atak,
743 J. Wouters, S. Aerts, SCENIC: single-cell regulatory network inference and
744 clustering, *Nat Methods* 14(11) (2017) 1083-1086.
- 745 [23] M.C. Oldham, S. Horvath, D.H. Geschwind, Conservation and evolution of gene
746 coexpression networks in human and chimpanzee brains, *Proceedings of the*
747 *National Academy of Sciences of the United States of America* 103(47) (2006)
748 17973-8.
- 749 [24] M. Schlichtkrull, T.N. Kipf, P. Bloem, R. van den Berg, I. Titov, M. Welling,
750 Modeling Relational Data with Graph Convolutional Networks, in: A. Gangemi, R.
751 Navigli, M.-E. Vidal, P. Hitzler, R. Troncy, L. Hollink, A. Tordai, M. Alam (Eds.) *The*
752 *Semantic Web*, Springer International Publishing, Cham, 2018, pp. 593-607.
- 753 [25] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, Y.J.a.e.-p. Bengio,
754 Graph Attention Networks, 2017, p. arXiv:1710.10903.
- 755 [26] T. Hoang, J. Wang, P. Boyd, F. Wang, C. Santiago, L. Jiang, S. Yoo, M. Lahne,
756 L.J. Todd, M. Jia, C. Saez, C. Keuthan, I. Palazzo, N. Squires, W.A. Campbell, F.
757 Rajaii, T. Parayil, V. Trinh, D.W. Kim, G. Wang, L.J. Campbell, J. Ash, A.J. Fischer,
758 D.R. Hyde, J. Qian, S. Blackshaw, Gene regulatory networks controlling vertebrate
759 retinal regeneration, *Science (New York, N.Y.)* 370(6519) (2020).
- 760 [27] V. Ravi, B. Venkatesh, The Divergent Genomes of Teleosts, *Annual review of*
761 *animal biosciences* 6 (2018) 47-68.
- 762 [28] S.M. Glasauer, S.C. Neuhauss, Whole-genome duplication in teleost fishes and
763 its evolutionary consequences, *Molecular genetics and genomics : MGG* 289(6)
764 (2014) 1045-60.
- 765 [29] B. Tasic, Z. Yao, L.T. Graybuck, K.A. Smith, T.N. Nguyen, D. Bertagnolli, J.
766 Goldy, E. Garren, M.N. Economo, S. Viswanathan, O. Penn, T. Bakken, V. Menon, J.
767 Miller, O. Fong, K.E. Hirokawa, K. Lathia, C. Rimorin, M. Tieu, R. Larsen, T. Casper,
768 E. Barkan, M. Kroll, S. Parry, N.V. Shapovalova, D. Hirschstein, J. Pendergraft, H.A.
769 Sullivan, T.K. Kim, A. Szafer, N. Dee, P. Groblewski, I. Wickersham, A. Cetin, J.A.
770 Harris, B.P. Levi, S.M. Sunkin, L. Madisen, T.L. Daigle, L. Looger, A. Bernard, J.
771 Phillips, E. Lein, M. Hawrylycz, K. Svoboda, A.R. Jones, C. Koch, H. Zeng, Shared
772 and distinct transcriptomic cell types across neocortical areas, *Nature* 563(7729)
773 (2018) 72-78.
- 774 [30] I. Korsunsky, N. Millard, J. Fan, K. Slowikowski, F. Zhang, K. Wei, Y. Baglaenko,
775 M. Brenner, P.R. Loh, S. Raychaudhuri, Fast, sensitive and accurate integration of
776 single-cell data with Harmony, *Nat Methods* 16(12) (2019) 1289-1296.
- 777 [31] L. McInnes, J. Healy, J.J.a.e.-p. Melville, UMAP: Uniform Manifold Approximation
778 and Projection for Dimension Reduction, 2018, p. arXiv:1802.03426.

- 779 [32] X. Zhang, Y. Lan, J. Xu, F. Quan, E. Zhao, C. Deng, T. Luo, L. Xu, G. Liao, M.
780 Yan, Y. Ping, F. Li, A. Shi, J. Bai, T. Zhao, X. Li, Y. Xiao, CellMarker: a manually
781 curated resource of cell markers in human and mouse, *Nucleic acids research*
782 47(D1) (2019) D721-d728.
- 783 [33] The Gene Ontology resource: enriching a GOld mine, *Nucleic acids research*
784 49(D1) (2021) D325-d334.
- 785 [34] M. Ashburner, C.A. Ball, J.A. Blake, D. Botstein, H. Butler, J.M. Cherry, A.P.
786 Davis, K. Dolinski, S.S. Dwight, J.T. Eppig, M.A. Harris, D.P. Hill, L. Issel-Tarver, A.
787 Kasarskis, S. Lewis, J.C. Matese, J.E. Richardson, M. Ringwald, G.M. Rubin, G.
788 Sherlock, Gene ontology: tool for the unification of biology. The Gene Ontology
789 Consortium, *Nature genetics* 25(1) (2000) 25-9.
- 790 [35] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the
791 Inception Architecture for Computer Vision, 2016 IEEE Conference on Computer
792 Vision and Pattern Recognition (CVPR), 2016, pp. 2818-2826.
- 793 [36] D.P. Kingma, J.J.C. Ba, Adam: A Method for Stochastic Optimization,
794 [abs/1412.6980](https://arxiv.org/abs/1412.6980) (2015).
- 795 [37] N.X. Vinh, J. Epps, J. Bailey, Information Theoretic Measures for Clusterings
796 Comparison: Variants, Properties, Normalization and Correction for Chance, 11
797 (2010) 2837–2854.
- 798 [38] H. Ahrens, Lancaster, H. O.: The Chi-squared Distribution. Wiley & Sons, Inc.,
799 New York 1969. X, 366 S., 140 s, 13(5) (1971) 363-364.
- 800 [39] A. Butler, P. Hoffman, P. Smibert, E. Papalexi, R. Satija, Integrating single-cell
801 transcriptomic data across different conditions, technologies, and species, *Nature*
802 *biotechnology* 36(5) (2018) 411-420.
- 803 [40] F.A. Wolf, P. Angerer, F.J. Theis, SCANPY: large-scale single-cell gene
804 expression data analysis, *Genome Biol* 19(1) (2018) 15.
- 805 [41] R.J. Kinsella, A. Kähäri, S. Haider, J. Zamora, G. Proctor, G. Spudich, J.
806 Almeida-King, D. Staines, P. Derwent, A. Kerhornou, P. Kersey, P. Flicek, Ensembl
807 BioMarts: a hub for data retrieval across taxonomic space, *Database : the journal of*
808 *biological databases and curation* 2011 (2011) bar030.
- 809 [42] V.A. Traag, L. Waltman, N.J. van Eck, From Louvain to Leiden: guaranteeing
810 well-connected communities, *Scientific reports* 9(1) (2019) 5233.
- 811 [43] A.D. Diehl, T.F. Meehan, Y.M. Bradford, M.H. Brush, W.M. Dahdul, D.S. Dougall,
812 Y. He, D. Osumi-Sutherland, A. Ruttenberg, S. Santivijai, C.E. Van Slyke, N.A.
813 Vasilevsky, M.A. Haendel, J.A. Blake, C.J. Mungall, The Cell Ontology 2016:
814 enhanced content, modularization, and ontology interoperability, *Journal of*
815 *biomedical semantics* 7(1) (2016) 44.
- 816 [44] V.Y. Kiselev, A. Yiu, M. Hemberg, scmap: projection of single-cell RNA-seq data
817 across data sets, *Nat Methods* 15(5) (2018) 359-362.
- 818 [45] M.D. Robinson, D.J. McCarthy, G.K. Smyth, edgeR: a Bioconductor package for
819 differential expression analysis of digital gene expression data, *Bioinformatics*
820 (Oxford, England) 26(1) (2010) 139-40.
- 821 [46] L. Haghverdi, A.T.L. Lun, M.D. Morgan, J.C. Marioni, Batch effects in single-cell
822 RNA-sequencing data are corrected by matching mutual nearest neighbors, *Nature*

823 biotechnology 36(5) (2018) 421-427.

824

825

826

827 **Figure 1. Overview of CAME.** (A) The architecture of the heterogeneous graph
828 neural network in CAME. The scRNA-seq data of both reference and query
829 species and their homology genes are encoded as a heterogeneous cell-gene
830 graph. The cell-gene edge indicates that the cell has non-zero expression of
831 the gene. The gene homologous mappings are represented by a gene-gene
832 bipartite graph with each edge indicating a gene homology. Note that the
833 homologous gene mappings can be many-to-many homologies. To preserve
834 the intrinsic data structure, the within-species cell-cell edges are adopted where
835 an edge between a pair of cells indicates that one is the k nearest neighbor of
836 the other (k=5 by default). The heterogeneous graph and the gene expression
837 profiles are input to CAME, passing through the inductive embedding layer, the
838 recurrent relational graph neural network, and the graph classifier with attention
839 mechanisms. The model is trained by minimizing the cross-entropy loss
840 calculated between the model prediction and the given labels of the reference
841 cells in an end-to-end manner. (B) Graph spatial convolutions for six different
842 types of edges including “a cell expresses a gene”, “a gene is expressed by a
843 cell”, “cell-cell similarity”, “gene-gene homology”, “cell self-loop” and “gene self-
844 loop” with the edge type-specific convolution weights. (C) Heterogeneous graph
845 attention classifier on the last layer, where each cell pays different attention to
846 its neighbor genes. The output cell-type probabilities are calculated by the
847 weighted sum of the neighbor-gene embeddings, followed by the softmax
848 normalization. The attention weights are calculated from the concatenated cell
849 and gene embeddings with a linear transformation, followed by activation and
850 the softmax normalization among the neighbor-genes of the cell. (D) The output
851 of CAME includes the probabilistic cell-type assignment of the query species,
852 as well as low-dimensional embeddings of the cells and genes from both
853 species. The gene embeddings are used for joint module extraction that allows
854 inter-species comparison of conservative or divergent characteristics.

855

856

857 **Figure 2. Benchmarking cross-species cell-type assignment performance**
858 **of CAME.** (A and B) Performance comparison of CAME and the six
859 benchmarking approaches in terms of cell-typing accuracy on 139 pairs of
860 cross-species scRNA-seq datasets (A) and on 510 pairs of cross-species
861 scRNA-seq datasets that associated with zebrafish, where each point
862 represents a pair of cross-species datasets and is colored by tissue. The
863 notation “X-Y” indicates that X is the reference and Y is the query. H: Human,
864 M: Mouse, C: Chick, Z: Zebrafish. (C) Performance comparison of the
865 classification accuracies of CAME and the six benchmarking methods on
866 different down-sampling rates (0.75, 0.5, 0.25, 0.1) for read counts.

867

868

869 **Figure 3. Alignment comparison of cell embeddings across datasets by**
870 **CAME and five benchmarking methods. (A)** The UMAP plots of the cell
871 embeddings by CAME and five benchmarking integration methods on the
872 scRNA-seq data from turtle (reference) and mouse (query) brains. Cells are
873 colored by their dataset identities (the first row) or cell type (the second row).
874 **(B)** Similar settings to (A). Here the reference datasets are the human
875 pancreatic scRNA-seq data from eight batches by five different platforms and
876 the query is from mouse pancreas cells. The UMAP plots of the third row
877 showed the reference datasets, colored by batch identities.

878

879 **Figure 4. Application of CAME to human and mouse scRNA-seq data of**
880 **brain cells. (A)** The predicted cell-type probabilities for each cell (each column)
881 in the mouse brain scRNA-seq data. The gene expressions of the human brain
882 were taken as the reference. Each row indicates a cell type in human data. OPC
883 is short for “oligodendrocyte precursor cells”, SMC is short for “smooth muscle
884 cell”, and VLMC is short for “vascular and leptomenigeal cell”. **(B)** The top
885 homologous DEG expressions of oligodendrocytes and (predicted) OPCs in
886 human and mouse data, including several marker genes reported by previous
887 literature (collected from CellMarker, colored by red or blue). **(C)** Cross-species
888 alignment of the gene embeddings output by CAME, where each dot represents
889 a gene and each edge indicates the homology between a pair of genes. Genes
890 shared between species are colored by light-blue (human) or pink (mouse)
891 while the other genes are colored by dark-blue (human) or dark-red (mouse).
892 **(D)** The UMAP plots of gene embeddings showing the average expression
893 patterns (z-scored across cell-types for each gene) of four cell types (excitatory
894 neurons, inhibitory neurons, oligodendrocytes, OPCs) of human and mouse
895 brains, where the color of each dot is scaled by the expression level of that cell
896 type in the gene. **(E)** Abstracted graph of the heterogenous cell-gene graph,
897 each node represents a cell type (pink) or a gene module (light blue). The size
898 of a node is scaled by the number of single cells in that type or the number of
899 genes in that gene module. The width of an edge is scaled by either the
900 normalized mean expression levels of a cell type in the connected gene module
901 or the conservancy of inter-species gene modules based on the gene
902 embeddings learned by CAME. **(F)** Gene modules detected by joint module
903 extraction of genes from humans (above) and mice (below).

904

905 **Figure 5. Application of CAME to human and macaque scRNA-seq data**
906 **during spermatogenesis. (A)** The predicted cell-type probabilities for each
907 macaque testicular cell (each column). The gene expressions of human testis
908 were taken as the reference. Each row indicates a cell type in the human data.
909 **(B)** The UMAP plots of cell embeddings output by CAME, colored by datasets
910 (left) or cell type (right). **(C)** 2D visualization of gene embeddings showing the
911 average expression patterns (z-scored across cell-types for each gene) of the
912 four stages across spermatogenesis, where each point represents a gene and
913 the color of each scatter is scaled by the expression level of that cell type in the
914 gene. **(D)** Abstracted graph of the heterogeneous cell-gene graph. Each node
915 represents a cell type (pink) or a gene module (light blue). The size of a node
916 is scaled by the number of single cells in that type or the number of genes in
917 that gene module. The width of an edge is scaled by either the normalized mean
918 expression levels of a cell type in the connected gene module or the
919 conservancy of inter-species gene modules based on the gene embeddings
920 learned by CAME. **(E)** Gene modules detected by joint module extraction of
921 genes from humans (above) and macaques (below).

922

923

Fig. 1

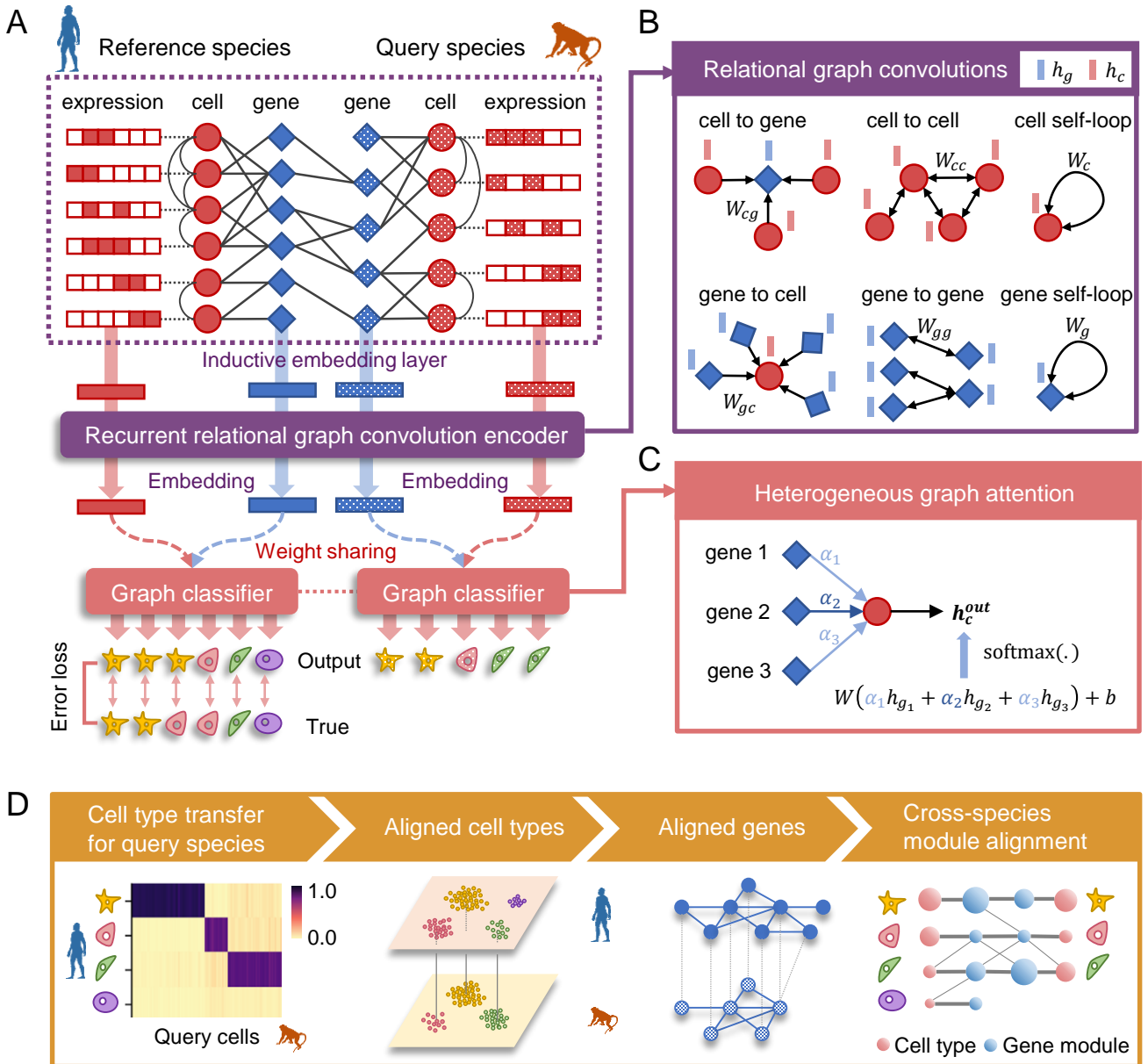


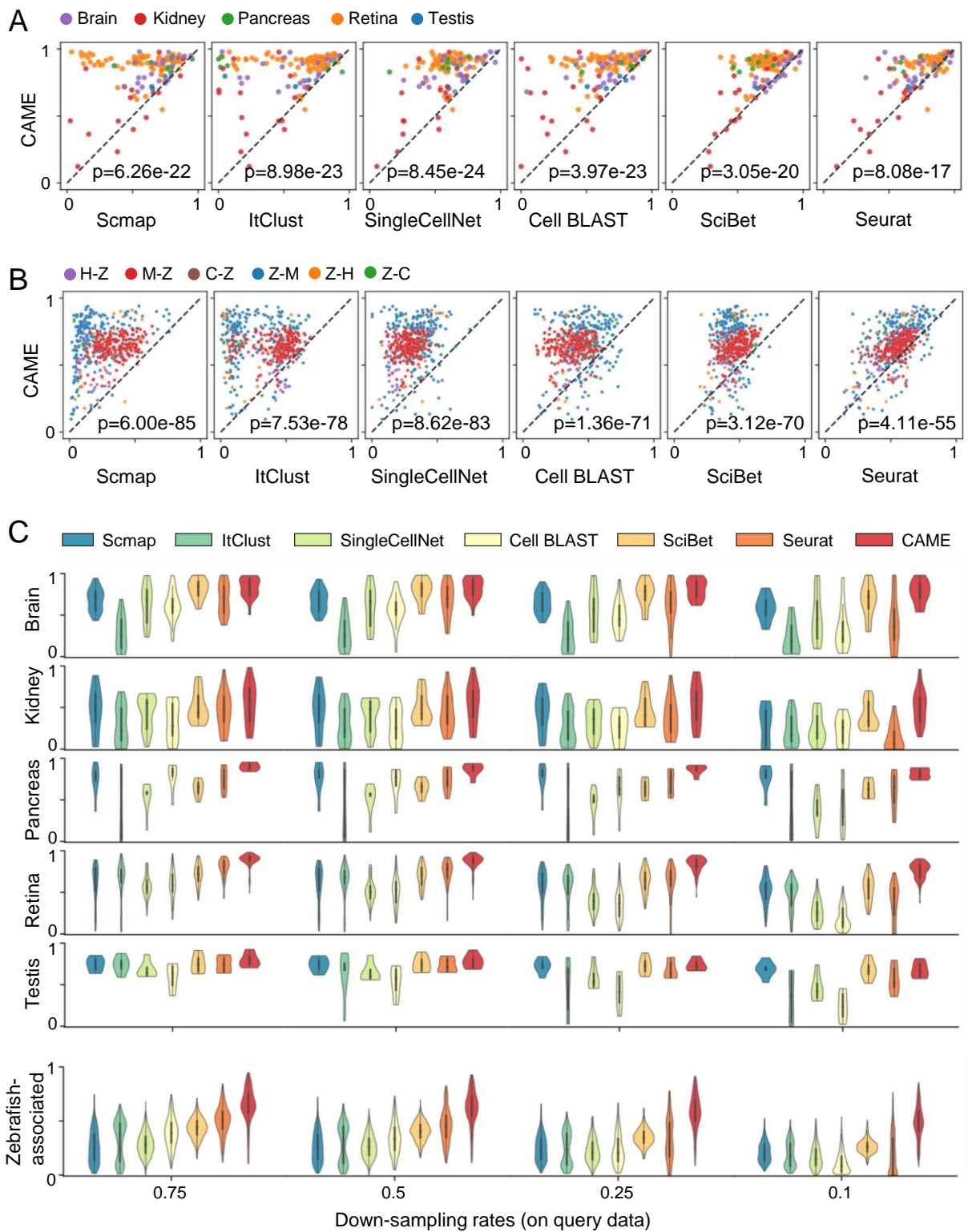
Fig. 2

Fig. 3

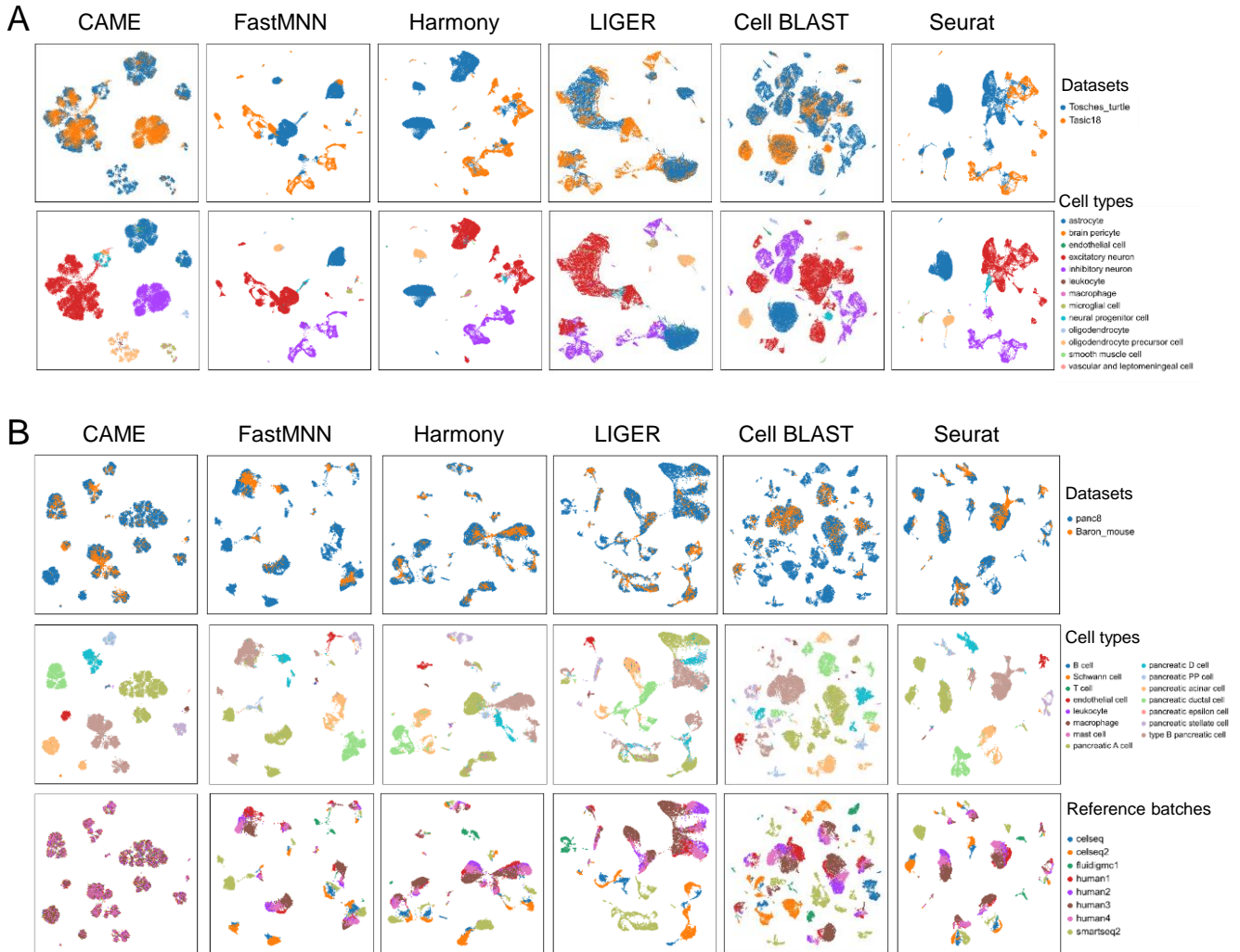


Fig. 4

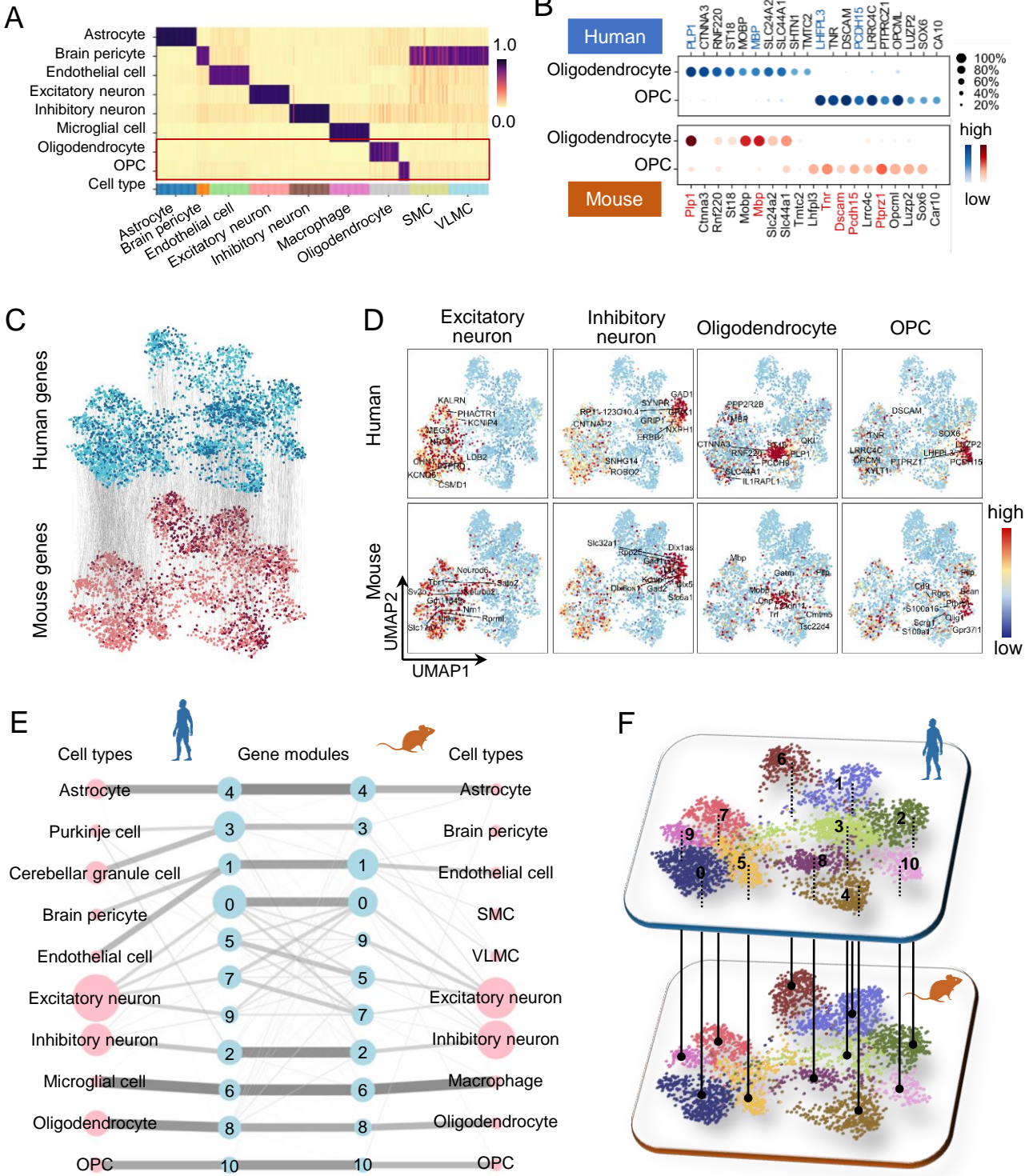


Fig. 5

