

Enhancement of network architecture alignment in comparative single-cell studies

Clemens Schächter^{*1,*}, Martin Treppner^{†1}, Maren Hackenberg^{‡1}, Hanne Raum^{§2},
Joschka Bödecker^{¶2,3}, and Harald Binder^{||1,4,5,*}

¹Institute of Medical Biometry and Statistics (IMBI), Faculty of Medicine and Medical
Center – University of Freiburg, Germany, ,

²Neurorobotics Lab, Dept. of Computer Science – University of Freiburg, Germany

³BrainLinks-BrainTools CRIION - Collaborative Research Institute Intelligent Oncology

⁴Freiburg Center for Data Analysis and Modelling – University of Freiburg, Germany

⁵CIBSS, Centre for Integrative Biological Signalling Studies – University of Freiburg,
Germany

*Corresponding author

October 6, 2024

Abstract

Animal data can provide meaningful context for human gene expression at the single-cell level. This can improve cell-type detection and clarify how well animal models represent human biology. To achieve this, we propose a deep learning approach that identifies a unified latent space to map complex patterns between datasets. The proposed method is tested to facilitate information transfer in liver, adipose tissue, and glioblastoma datasets from various animal models. Our results are robust for small datasets and large differences in the observed gene sets. Thus, we reliably uncover and exploit similarities between species to provide context for human single-cell data.

*clemens.schaechter@uniklinik-freiburg.de

†Martin-Treppner@gmx.de

‡maren.hackenberg@uniklinik-freiburg.de

§raum@informatik.uni-freiburg.de

¶j.boedeck@informatik.uni-freiburg.de

||harald.binder@uniklinik-freiburg.de

22 **Keywords:** Cross-species alignment, Model organisms, Deep learning, Transfer learning, Variational
23 autoencoder, Single-cell RNA sequencing, Comparative genomics.

24 1 Background

25 Model organisms are crucial in advancing biomedical research by offering advantages such as easy ge-
26 netic manipulation and access to datasets from a variety of experimental contexts [1]. As a popular choice,
27 mouse models have significantly contributed to the study of human diseases [2], including diabetes [3],
28 glioblastoma [4], and non-alcoholic fatty liver disease [5]. However, translating experimental findings to
29 humans is challenging owing to biological differences between species. Efforts to bridge this evolutionary
30 gap include engineered mouse models that replicate human biology more closely [6]. The emergence of
31 single-cell RNA sequencing (scRNA-seq) has also opened up opportunities for deep learning approaches
32 to compare experimental findings across species.

33 Transfer learning techniques have established themselves as powerful tools for sharing information be-
34 tween scRNA-seq datasets. These approaches often use encoder-decoder architectures to compress
35 datasets into a low-dimensional manifold. Examples include Cell BLAST [7] and ItClust [8], which anno-
36 tate and cluster cells based on knowledge transfer from reference datasets.

37 Architecture surgery techniques adjust network architectures according to the characteristics of different
38 datasets. After pretraining, additional neurons are inserted into the encoder and decoder input layers.
39 These neurons correct for unseen batch effects in the new data, while all other weights remain fixed
40 during subsequent training. This approach, pioneered by scArches [9], now spans a diverse set of mod-
41 els [10–12]. Despite the method’s success, two primary challenges remain unaddressed for datasets of
42 different species (Figure 4).

43 First, some genes lack orthologs in other genomes, which requires different interpretations of certain in-
44 put nodes in their neural network architectures. For example, 20% of human protein-coding genes and
45 a significant percentage of small and long noncoding RNAs lack one-to-one mouse orthologs [13]. To
46 enable training, architecture surgery-based approaches restrict datasets to orthologous genes or zero-fill
47 missing values. Outside of architecture surgery, some models like SATURN [14] and TACTiCS [15] match
48 genes via protein sequences with transformer-based language models.

49 The second challenge is that biological similarities between cells do not always translate into similar gene
50 expression patterns, which can vary significantly between species [13]. Therefore, neural networks may
51 struggle to recognize similar cells.

52 To account for differences between gene sets and expression levels, we introduce scSpecies. Our ap-
53 proach pretrains a conditional variational autoencoder-based model [16] and fully reinitializes the encoder
54 input layers and the decoder network during fine-tuning. Architecture alignment is guided by a nearest

neighbor search performed on homologous genes, which estimates the similarity between cells in both datasets. This incentivizes our model to map biologically related cells into similar regions of the latent space. The neighbor search requires only a small subset of observed genes to be homologs, while all remaining genes can have no relationship at all. Moreover, scSpecies enables nuanced comparisons of gene expression profiles by generating gene expression values for both species from a single latent variable.

We tested our method on data from various species and organs, including liver cells [17], white adipose tissue cells [18], and glioblastoma immune response cells [19]. Our results demonstrate that scSpecies effectively aligns network architectures and latent representations. We improve upon cell label transfer from the initial nearest neighbor search and existing architecture surgery approaches when measured in terms of accuracy and multiple clustering metrics.

2 Results

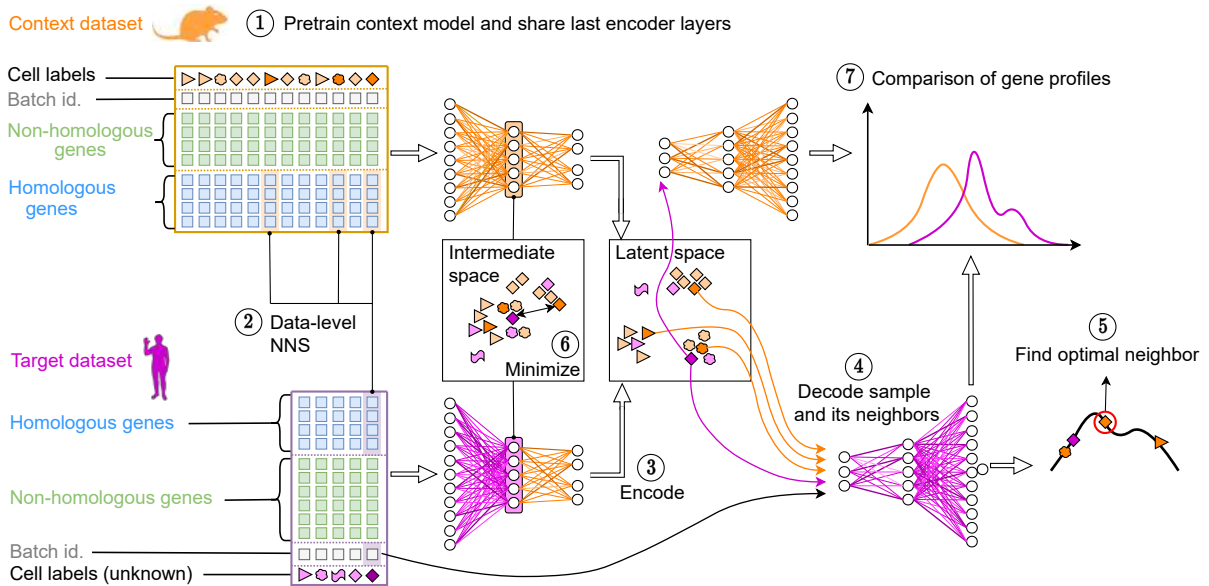


Figure 1: Graphical representation of the scSpecies workflow. Step 1: The encoder and decoder neural networks are trained on the dataset of the context species. The weights of the last encoder layers are incorporated into the encoder model for the target species. Step 2: A nearest neighbor search is performed on the shared genes of the context and target dataset. This identifies a set of k context neighbors for every target cell. Step 3: The cells of the target dataset are encoded into the latent space. For cells with high agreement among the cell labels of their neighbors, we retrieve the latent variables of their neighbors. Step 4: The latent values of their k neighbors are passed to the decoder together with the human batch label. Step 5: The optimal candidate among the k neighbors is chosen as the cell with the highest log-likelihood. Step 6: The distance between the optimal candidate and the intermediate representation of its target cell is minimized. Step 7: After training, normalized gene expression profiles can be compared by decoding latent variables with both decoder networks. Additionally, labels can be transferred via the aligned latent representation.

We present scSpecies, a tool for researchers who wish to use one scRNA-seq dataset as a context for

68 another from a different species. In the following, the dataset of the model organism is referred to as the
69 'context dataset', and the dataset of the target organism is referred to as the 'target dataset'. scSpecies
70 aligns context scRNA-seq datasets with human target data, enabling the analysis of similarities and dif-
71 ferences between the datasets.

72 In addition to the context and target datasets, the model requires a sequence containing indices of ho-
73 mologous genes, indicator variables for batch effects, and cell type labels for the context dataset.

74 The proposed workflow (Figure 1) aligns the network architectures of two single-cell variational infer-
75 ence (scVI) [20] models in a pretraining strategy. In scVI, encoder neural networks map gene expression
76 vectors into a compressed latent space separating cells by biological features. Conversely, a decoder
77 maps from this low-dimensional representation onto parameters of a negative binomial distribution to (re-
78)generate gene expression data.

79 First, our proposed approach pretrains a scVI model on the context dataset. Afterwards, the last encoder
80 layers are transferred into a second scVI model for the target species. The aim of this architecture transfer
81 is to share learned information within the network weights between datasets and species. During subse-
82 quent fine-tuning, the shared weights remain frozen while all other weights are optimized.

83 Unlike existing architecture surgery approaches, we align the architectures in a reduced intermediate fea-
84 ture space instead of at the data level. This approach is inspired by the notion of midlevel features from
85 computer vision [21, 22]. These represent abstractions of the input image learned by neural networks
86 in their intermediate layers. Midlevel features combine individual elements into more general structures,
87 such as contours, specific shapes, or parts of objects. Transfer learning approaches then retrain the
88 last layers to transition these intermediate representations into task-specific network outputs for different
89 datasets [23].

90 Unlike images, scRNA-seq datasets lack ordered patterns as gene expression vectors can be permuted
91 without changing their information content. Nevertheless, the first encoder layers translate dataset-
92 specific features, such as influences of experimental batches or interactions between observed genes,
93 into a higher abstraction level (Figure 5). The resulting representation may correspond to more funda-
94 mental cell properties that are less perceptible to noise and systematic differences between species.

95 To connect the new encoder layers with the pretrained structure, we identify sets of similar cells through a
96 nearest neighbor search performed on homologous genes. Afterward, scSpecies minimizes the distance
97 between a target cell's midlevel representation and a suitable candidate from its set of neighbors. The
98 model determines the most suitable context cell as the candidate whose decoded latent representation
99 yields the highest log-density value at the location of the target cell within the decoder's distribution. To
100 counter misclassifications, we align midlevel features for only those target cells whose context neighbors
101 have high agreement in their cell labels.

102 During model fitting, we thus encode similarity information both at the original data level and at the level

103 of learned features. The aligned latent space then captures cross-species similarity relationships based
 104 on the fitted model, which facilitates information transfer across species.

105 2.1 scSpecies aligns architectures across species

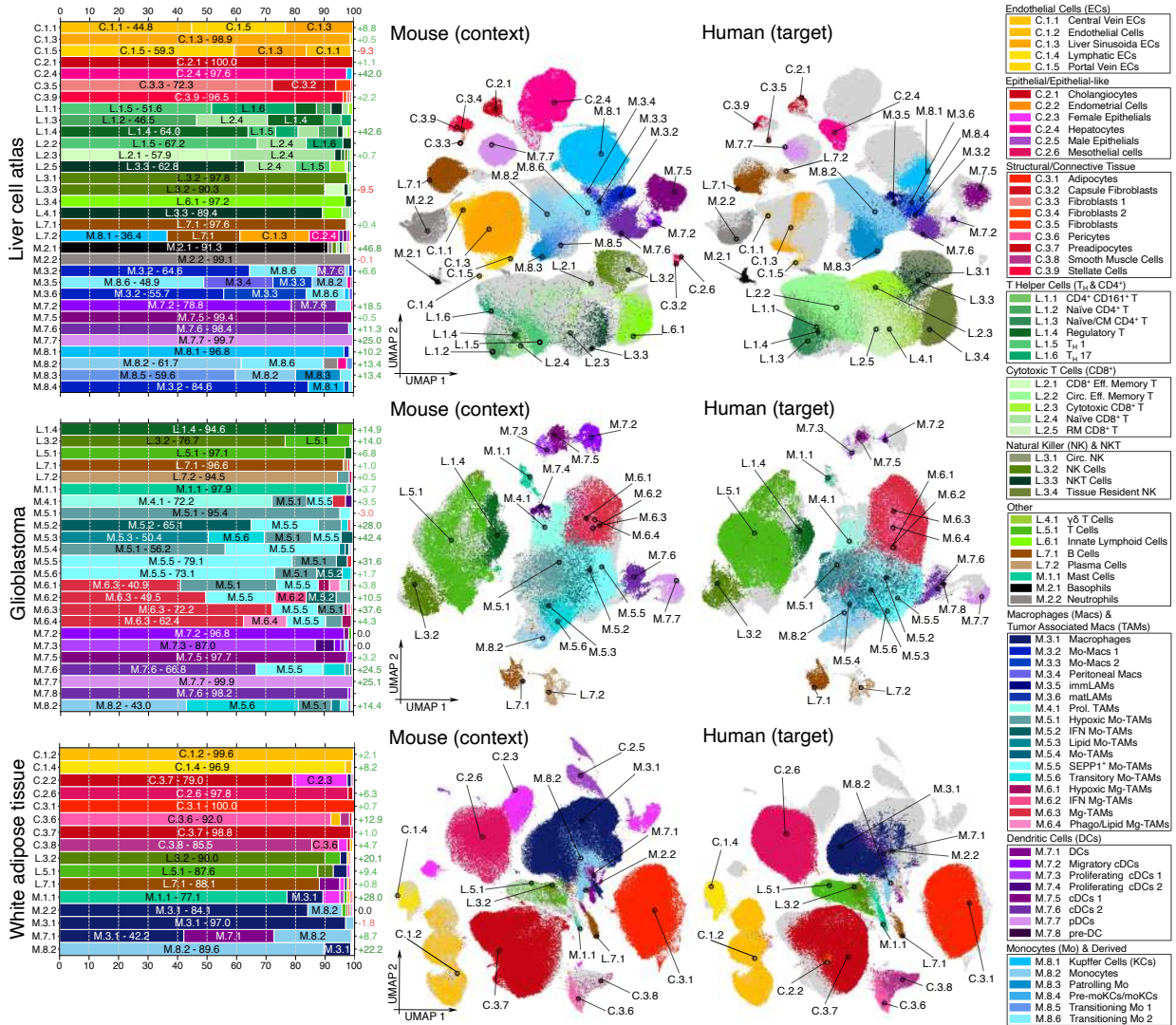


Figure 2: Visualization of the aligned representations for three dataset pairs obtained by training scSpecies with a set of 25 neighbors. We color cells by fine cell type labels for the liver and glioblastoma datasets, and by coarse cell labels for the adipose tissue dataset. On the left, the bar plots indicate the accuracy of cell label transfer through a nearest neighbor search in the aligned latent space. The left y-axis labels indicate cell type codes corresponding to human cell labels. These codes are referenced in the legend. The bars contain the frequency of assigned mouse cell labels. The results are averaged over five random seeds. The left y-axis labels indicate improvement in accuracy for shared cell types over the data-level nearest neighbor search. In addition to the bar plots, the UMAP coordinates of the aligned latent representations are visualized. The lymphoid cell types are colored in green and brown; the myeloid cell types are colored blue and purple; and the CD45⁻ cell types are colored red, pink and yellow. The cells from the other dataset are indicated in a light gray.

106 We applied the scSpecies workflow to three mouse-human dataset pairs containing liver cells, white
 107 adipose tissue cells, and immune response cells to glioblastoma.

108 We visually examined alignment through UMAP coordinates [24] of the combined latent variables of

109 dataset pairs (Figure 2). The 2D representation showed biologically meaningful alignment of the cells.
110 Cell types without context counterparts aligned with related cell types or formed distinct clusters.
111 To facilitate label and information transfer for target cells, we conducted a second nearest neighbor search
112 on the shared latent representation of both datasets. Afterwards, we inferred target cell labels from their
113 set of latent context neighbors via majority voting. For labels at the subcell type resolution, the accuracy
114 was 73% for liver, 49% for adipose tissue, and 69% for glioblastoma datasets. Misclassifications mostly
115 occurred within biologically related cells belonging to the same overarching cell type. For broader cell type
116 labels, accuracy increased to 92% for the liver, 82% for the adipose tissue, and 80% for the glioblastoma
117 dataset. These values represent significant improvements upon the data-level nearest neighbor search
118 and existing architecture surgery approaches (Table 7). We also calculated the adjusted Rand index and
119 adjusted mutual information and observed improvements in these metrics.
120 We observed a greater increase in label transfer accuracy for cell types with noisy data-level nearest
121 neighbor search but clear separation in their pretrained latent space. For example, the initial neighbor
122 search matched less than half of all human liver basophils (cluster M.2.1) with mouse counterparts. This
123 value improved to over 90% through our method. However, in the adipose tissue datasets, neither the con-
124 text scVI model nor the nearest neighbor search separated dendritic cells, monocytes, and macrophages.
125 Thus, scSpecies could not separate these cell types either.
126 The results were consistent over architecture variations and averaged over five random seeds; however,
127 for cell types with noisy neighbor search results, like hepatocytes or portal vein endothelial cells, misclas-
128 sifications of the whole cell type occurred in one random seed.
129 We also tested scSpecies in a scenario where the target dataset was small but equally diverse in terms of
130 cell types and batch effects. Specifically, we randomly sampled 5000 cells from the human liver dataset
131 and trained the model to align with the full mouse context dataset. We repeated sampling and training ten
132 times and obtained accuracy scores of 88% and 68% for coarse and fine cell labels, respectively, which
133 still indicates reasonable performance.

134 **2.2 The nearest neighbor search is an important component of scSpecies.**

135 We explored the importance of incorporating the nearest neighbor search into scSpecies. (Table 7) With-
136 out this component, we observed misaligned latent representations and significantly reduced label trans-
137 fer accuracy. Initializing the inner encoder layers with random, frozen weights yielded similar results to
138 using the pretrained structure. This implies that without an explicit neighbor alignment component, trans-
139 ferred layers were treated like random nuisances.

140 Training with one neighbor forced the model to align some cells with mismatched counterparts as the
141 approach could not choose from a set of suitable options. We observed meaningful alignment but with
142 reduced performance.

143 Training with 25 neighbors improved the results noticeably on all datasets. To investigate the preferred
144 candidate choice, we tracked the cell prototypes during alignment. We created context and target pro-
145 totype cells consisting of empirical median gene expression values within a cell type. For each target
146 prototype, we included all context prototypes within its set of candidates and tracked their log-likelihoods
147 during alignment (Figure 10). At onset, the likelihoods for all prototypes were nearly equal. This resulted
148 in alignment driven by chance favoring cell candidates of the most occurring cell label. For cell types
149 with a noisy neighbor set, corrections during later training stages eventually aligned them with appropri-
150 ate prototypes. We observed this with hepatocytes, migratory cDCs, and basophils, which had nearest
151 neighbor search accuracies of 56%, 61%, and 45%, respectively. The cell types where the neighbor
152 search yielded predominantly incorrect results did not align correctly, such as killer T cells and cytotoxic
153 CD8⁺ cells, which had initial accuracies of only 11% and 1%, respectively.
154 Finally, alignment with a large neighbor set caused neglect of rare cell types, resulting in lower corre-
155 sponding accuracy scores. Metrics such as the adjusted Rand index and adjusted mutual information
156 were comparable or improved, as they do not reflect different cell type label sizes.

157 **2.3 scSpecies can help to better separate latent cell clusters.**

158 To investigate the intermediate representations, we compared the clustering quality of intermediate repre-
159 sentations in unaligned and aligned scVI architectures. We found that clustering based on experimental
160 batches became increasingly mixed as the data progressed toward the latent space. In the unaligned
161 architectures, the Davies-Bouldin index (DBI) increased from 10 to 21.9 in the mouse context, and from
162 15.8 to 33.5 in the human liver dataset. Conversely, cell type clusters showed increasingly better separa-
163 tion, resulting in a DBI reduction from 4.6 to 1.6 and from 4.9 to 2.4 for the mouse and human datasets,
164 respectively (Figures 5,6,7).

165 This phenomenon is caused by the design of scVI, which removes batch influences to enforce a normal
166 distribution in the latent space. Batch patterns are added by the decoder through their provided labels.
167 However, scVI must separate cell types to reconstruct cell characteristics from the latent representation.
168 Yet, certain cell types in the human liver dataset, such as hepatocytes, stellate cells, and fibroblasts, are
169 predominantly associated with a single batch label. Consequently, the model inferred cell type information
170 from batch labels, removing biological characteristics from their latent variables. However, these cell types
171 were still separated in the intermediate spaces which are not regularized to follow a normal distribution.
172 Alignment adjusted the target encoder architecture to the well-separated latent mouse context represen-
173 tation. This improved latent cell cluster separation, as measured by a decrease in DBI from 2.4 to 1.8.
174 For white adipose tissue and glioblastoma dataset pairs, clustering improvement was marginal, with a
175 decrease in DBI from 1.7 to 1.6 and from 2.2 to 2, respectively.

176 We also studied the effectiveness of directly aligning latent representations. Direct latent alignment does

177 not require access to the context model weights. However, we observed a decline in performance met-
 178 rics across all datasets. This underlines the potential of better alignment within the more information-rich
 179 midlevel feature spaces.

180 **2.4 scSpecies can align datasets of multiple species.**

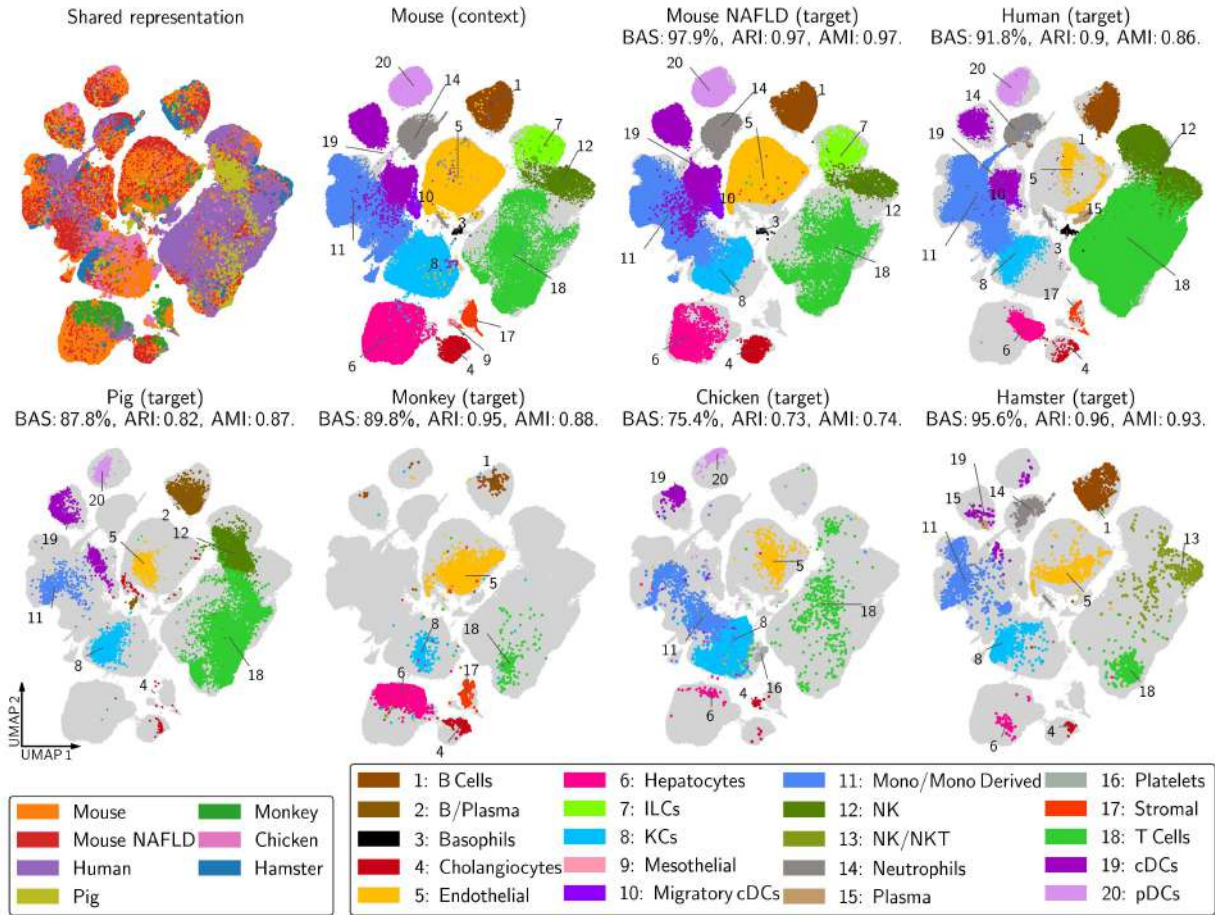


Figure 3: We utilized scSpecies to obtain an aligned liver cell landscape that spans multiple species. The mouse dataset serves as a context for each species.

181 We employed scSpecies to simultaneously align liver cells from mice with fatty liver disease, humans,
 182 pigs, monkeys, chickens, and hamsters, using a context dataset of healthy mice (Figure 3).

183 We successfully obtained aligned latent representations across species, despite fewer than half of the
 184 genes having mouse orthologs in some datasets.

185 An intriguing application of scSpecies is the potential to align datasets with very limited gene coverage, or
 186 even when there is no overlap in the observed gene set. This can be achieved by aligning each dataset
 187 to a comprehensive context dataset that shares a common gene set with both.

188 However, a limitation of this approach is its inability to align cell types not present in the context dataset.
 189 For example, plasma cells, which were absent from the mouse dataset, were not aligned across the
 190 human, pig, and hamster datasets.

2.5 scSpecies offers insights into the genetic manifestations of cells across species.

To better understand the similarities and differences between context and target datasets, e.g., to clarify in what aspects an animal might be a good model of human biological processes, we extended our analysis from the latent space to the data level. Here, we compared the reconstructed gene expression profiles and assigned relevance scores to the input genes.

We decoded latent representations using both decoder models to obtain normalized gene expression vectors for each species. These vectors allow us to compare and analyze the gene expression profiles of cells that have similar underlying biological properties. This analysis benefits from the correspondence between latent representations of both species, which is difficult to establish at the data level.

For our investigation, we focused on cell types present in both the mouse and human liver datasets. We assessed Log2Fold changes (LFCs) in normalized gene expression vectors, which indicate differences in gene expression levels between species. We also calculated the probability of observing genes as differentially expressed when sampling from the latent distribution of a cell type (Figure 8). Averaging across cell types revealed that 56% of the genes exhibited an LFC value above one. Among these, 15% of mouse genes were upregulated and 21% were downregulated compared with their human counterparts in over 90% of decoded cells. With an LFC threshold of two, 24% of genes had an LFC outside this boundary. With an LFC value of 0.4, a substantial 82% of genes showed an LFC outside this boundary. These results agree in magnitude with [25], who found an LFC value of greater than 0.4 in 78% of genes comparing humans with non-alcoholic liver disease and mice on a high-fat diet.

For white adipose tissue datasets, 50%, and for glioblastoma datasets, 47% of genes exhibited an LFC value greater than one.

We compared this with training on context-target dataset pairs of healthy mice and mice with liver disease. Here, only 22% of genes had an LFC value above one. Of those differentially expressed genes, 4% and 5% were upregulated and downregulated in more than 90% of samples. Only 6% of genes had an LFC over two, while 55% of genes showed LFC values above 0.4.

We extended our study by calculating relevance scores via Layer-wise relevance propagation (LRP) [26] (Figure 9). These scores measure each gene's contribution to a cell's latent value, offering insights into the learned significance of specific genes across different cell types and species. LRP was recently used to explain neural network predictions on scRNA-seq data [27].

First, we found no significant difference in relevance scores between non-homologous and shared genes, suggesting that training networks on a reduced gene set omits informative parts of the data.

Second, we found that the relevance scores were correlated with the gene expression levels. For the mice and human liver datasets, we found a Spearman's ρ between the expression level of genes and their relevance scores of 0.67 and 0.69 and a Pearson correlation coefficient of 0.63 and 0.71. This

226 suggests that differences in gene expression translate into relevant features for the neural networks. A
227 gene with high relevance scores across most cell types was *MALAT1*, which is highly conserved across
228 mammals [28].

229 **3 Discussion**

230 We introduced scSpecies, a novel deep learning approach designed to align neural network architectures
231 across different species. Aligning such architectures has been a challenging task due to differences in
232 genomes between species and variations in gene expression levels, even among homologous genes. Key
233 features of scSpecies include the retraining of the first encoder layers and integrating a nearest neighbor
234 search within the model. By focusing on the alignment of intermediate neural network layers rather
235 than the input layers, scSpecies captures more abstract biological properties that are less affected by
236 noise and species-specific variations. Additionally, the integration of a nearest neighbor search based on
237 homologous genes leverages model-based similarity information to guide the alignment process, ensuring
238 that biologically similar cells are mapped closely in the latent space.

239 Our results demonstrate that scSpecies effectively aligns scRNA-seq data from diverse species, including
240 mouse, human, pig, monkey, chicken, and hamster, across various tissues such as liver, white adipose
241 tissue, and glioblastoma cells. The method shows robust performance even when the datasets have a
242 limited number of shared genes or when the target dataset is small but diverse.

243 However, one limitation of the presented method is that cell types unique to the target dataset tend to be
244 aligned with biologically close cell types in the context dataset instead of being identified as new clusters
245 by the model. This could lead to misinterpretation of species-specific cell populations. Additionally, when
246 creating a collection of multiple species, cell types not present in the context dataset will not align across
247 species that exhibit them. To avoid misalignment, the context dataset should therefore encompass all
248 suspected cell types of the reference datasets.

249 There remain multiple potential directions for further development of our approach. While we initially
250 tested scSpecies with a scVI base model, the method could be easily adapted to other CVAE-based
251 models in the future. Furthermore, scSpecies could be extended to handle multimodal datasets, such
252 as those integrating scRNA-seq with protein expression data (CITE-seq). Our method would also benefit
253 from a direct metric that identifies cell types unique to the target datasets and detects cells that may be
254 misclassified due to noisy nearest neighbor search results.

4 Conclusions

We have introduced scSpecies, a novel deep learning approach that extends architecture surgery techniques to align scRNA-seq datasets across species. By retraining the first encoder layers, our method overcomes challenges posed by non-orthologous genes and divergent gene expression patterns, enabling more accurate cross-species comparisons. By aligning datasets from multiple species — even with minimal gene overlap — scSpecies provides a framework to better understand and compare the cellular and molecular similarities and differences of scRNA-seq datasets across species. Therefore, we envision that our method could lead to more effective translation of experimental findings from model organisms to humans, ultimately advancing our understanding of human biology.

5 Methods

In the following, we represent multidimensional vectors using bold italics and scalar values in regular italics. Dataset elements are indicated with superscript indices, and vector positions with subscript indices. The context dataset is indicated by the subscript C and the target dataset by the subscript T . Superscripts and subscripts are omitted when they are exchangeable. Random variables are expressed in a sans-serif mathematical font, as in X, Z, L . We represent distributions of random variables with uppercase letters, such as P_Z , and their probability density functions with lowercase letters, like $p_Z(z)$. Conditional distributions are denoted as $P_{X|s} := P_{X|S=s}$. In the following, we briefly describe the scVI model, which we subsequently use as a core of our proposed approach.

5.1 Single cell variational inference

Consider a dataset $\mathbb{D} = \{(\mathbf{x}^{(i)}, \mathbf{s}^{(i)})\}_{i=1}^M$ obtained through a single-cell RNA sequencing experiment. The mathematical model behind scVI [20] assumes that gene expression count vectors \mathbf{x} , and batch indicator variables \mathbf{s} , correspond to observations of random variables X and S . The gene expression data distribution $P_{X|s}$ is conditioned on its batch effect $S = s$. This accounts for technical artifacts during data collection. Within an experimental batch, gene expression vectors are independent and identically distributed samples from $P_{X|s}$.

scVI models the data distribution within a parametric family. Building on conditional variational autoencoders [16], a latent variable model is introduced. The random variable Z , corresponding to the representation of a cell in the latent space \mathbb{R}^d , is employed to capture biological variability among cells in the dataset. The one-dimensional random variable L with latent space $\mathbb{R}_{>0}$ accounts for technical variability due to different library sizes. Within the model, data is generated by drawing samples for Z and L from a prior distribution $P_{Z,L|s}$. Then, gene expression data is generated by drawing from the sampling distribu-

286 tion $P_{X|z,l,s}$.

287 The data p.d.f. $p_{X|s}$ can be expressed by integrating the joint probability across the latent spaces and then
 288 applying the general product rule of probability,

$$p_{X|s}(\mathbf{x}) = \int_{\mathbf{z}} \int_l p_{X|z,l,s}(\mathbf{x}) p_{Z,L|s}(\mathbf{z}, l) d\mathbf{z} dl. \quad (1)$$

289 To approximate this integral, scVI performs variational inference on the intractable posterior distribution
 290 $P_{Z,L|x,s}$. Therefore, the posterior probability is approximated by a variational distribution, denoted as
 291 $Q_{Z,L|x,s} \approx P_{Z,L|x,s}$. Further, scVI applies a mean field approximation, where p.d.fs of both variational and
 292 prior distribution are factorized,

$$q_{Z,L|x,s}(\mathbf{z}, l) = q_{Z|x,s}(\mathbf{z}) q_{L|x,s}(l), \quad p_{Z,L|s}(\mathbf{z}, l) = p_Z(\mathbf{z}) p_{L|s}(l). \quad (2)$$

293 The prior P_Z is assumed to be independent of S and fixed as standard normal distribution $P_Z = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

294 The prior $P_{L|s}$ is set as a log-normal distribution $P_{L|s} = \text{LogNormal}(\mathbf{l}_\mu^\top \mathbf{s}, \mathbf{l}_{\sigma^2}^\top \mathbf{s})$. The prior parameters
 295 are derived from empirical batch means and variances of the observed log-library sizes. The variational
 296 distribution $Q_{Z|x,s}$ is chosen as a normal distribution $\mathcal{N}(\boldsymbol{\mu}_Z, \boldsymbol{\sigma}_Z^2 \mathbf{I}_d)$, and $Q_{L|x,s}$ is set as a log-normal
 297 distribution $\text{LogNormal}(\mu_L, \sigma_L^2)$.

298 The parameters for these distributions are determined by two encoder neural networks,

$$f_{\text{enc}Z}(\mathbf{x}, \mathbf{s}) = (\boldsymbol{\mu}_Z, \boldsymbol{\sigma}_Z) \text{ and } f_{\text{enc}L}(\mathbf{x}, \mathbf{s}) = (\mu_L, \sigma_L). \quad (3)$$

299 scVI obtains latent variables by sampling from the variational distributions through the reparametrization
 300 trick [29].

301 The sampling distribution $P_{X|z,l,s}$ for generating gene-expression data from a given latent variable is as-
 302 sumed to follow a Gamma-Poisson mixture, resulting in a negative binomial distribution. The correspond-
 303 ing decoder network outputs a denoised gene expression vector that sums to one.

$$f_{\text{dec}}(\mathbf{z}, \mathbf{s}) = \boldsymbol{\rho}, \quad \sum_{g=1}^N \rho_g = 1. \quad (4)$$

304 The value ρ_g provides an estimate of the percentage of transcripts in a cell that originate from gene g .
 305 Gene expression values x_g can be drawn from a negative binomial distribution $\text{NB}(l\rho_g, \theta_{g,s})$ parameter-
 306 ized by mean $l\rho_g$ and dispersion $\theta_{g,s}$. The dispersion parameter is constant for every gene across cells
 307 of batch s . To address the potential issue of dropout, a zero-inflated negative binomial distribution can
 308 be used to model count data. The dropout probability parameter π is also obtained from the decoder
 309 network. The weights of the three neural networks and the parameters $\theta_{g,s}$ are optimized simultaneously

310 by empirically estimating and maximizing the ELBO function

$$\text{ELBO}(\mathbf{x}, \mathbf{s}, \beta) = \mathbb{E}_{\mathbf{q}_{\mathbf{z}, \mathbf{l}}|\mathbf{x}, \mathbf{s}} [\log p_{\mathbf{x}|\mathbf{z}, \mathbf{l}, \mathbf{s}}(\mathbf{x})] - \beta (D_{\text{KL}}[\mathbf{Q}_{\mathbf{z}|\mathbf{x}, \mathbf{s}} \parallel \mathbf{P}_{\mathbf{z}}] + D_{\text{KL}}[\mathbf{Q}_{\mathbf{l}|\mathbf{x}, \mathbf{s}} \parallel \mathbf{P}_{\mathbf{l}|\mathbf{s}}]) \quad (5)$$

311 on mini batches $\mathbb{M} \subset \mathbb{D}$.

312 5.2 The scSpecies approach

313 We consider a scenario involving two scRNA-seq datasets,

$$\mathbb{D}_C = \left\{ (\mathbf{x}_C^{(i)}, \mathbf{s}_C^{(i)}, c_C^{(i)}) \right\}_{i=1}^{M_C} \text{ and } \mathbb{D}_T = \left\{ (\mathbf{x}_T^{(j)}, \mathbf{s}_T^{(j)}) \right\}_{j=1}^{M_T}. \quad (6)$$

314 Their data points consist of gene expression measurements \mathbf{x} and batch indicator variables \mathbf{s} from a
 315 context species C and a target species T . Furthermore, context count vectors are clustered into distinct
 316 groups based on cell type labels c_C , whereas target labels c_T are unknown.

317 The count vectors from both datasets share a gene subset \mathbf{h} comprising count values from homologous
 318 genes,

$$\mathbf{x} = \underbrace{(x_1, \dots, x_H)}_{\mathbf{h} \text{ homologous}} \underbrace{(x_{H+1}, \dots, x_N)}_{\text{non-homologous}}. \quad (7)$$

319 The number of non-homologous genes can differ in both datasets, either because a gene has no ortholog
 320 in the genome of the other species or because it is not observed within the dataset. Therefore, gene
 321 expression vectors can be of different dimension, $N_C \neq N_T$.

322 To map both datasets into a unified latent space, we define separate scVI models for each dataset,

$$\text{scVI}^C = (f_{\text{encZ}}^C, f_{\text{encL}}^C, f_{\text{dec}}^C), \quad \text{scVI}^T = (f_{\text{encZ}}^T, f_{\text{encL}}^T, f_{\text{dec}}^T). \quad (8)$$

323 We divide the training procedure for scSpecies into three steps: Training of the context scVI model,
 324 followed by an initial data-level nearest neighbor search, and alignment of context and target latent rep-
 325 resentations.

326 5.2.1 Pretraining on the context dataset

327 First, the model scVI^C is trained on the context dataset by minimizing its negative ELBO function. Follow-
 328 ing training, the architecture of the encoder network for the latent variable \mathbf{z} is split up into two parts:

$$f_{\text{encZ}}^C = f_{\text{outer}}^C \circ f_{\text{inner}}^C. \quad (9)$$

329 The outer part f_{outer}^C consists of the first L layer functions and maps data from the input space \mathcal{X}_C to
 330 an intermediate feature space \mathcal{T} . The inner part, f_{inner}^C , consists of the last M layers. It encodes an
 331 intermediate representation onto the variational parameters with subsequent reparametrization into the
 332 latent space \mathcal{Z} . We incorporate this inner encoder part into the encoder architecture of scVI^T ,

$$f_{\text{encZ}}^T = f_{\text{outer}}^C \circ f_{\text{inner}}^T. \quad (10)$$

333 5.2.2 Nearest neighbor search

334 When the first layers are initialized randomly, the target model scVI^T cannot leverage the learned structure
 335 in its subsequent encoder layers. To leverage the learned weights, we incentivize alignment of interme-
 336 diate target representations with intermediate features of similar context cells. This leads to an aligned
 337 latent space as layer weights mapping from the intermediate space to the latent space are not updated.
 338 To quantify similarity and establish a direct correspondence between cells of context and target dataset,
 339 we perform a nearest neighbor search on the shared homologous gene subset \mathbf{h} . The nearest neighbors
 340 serve as a set of candidates for every target cell from which the model can choose a best fit to align their
 341 intermediate representations during the last training phase.

342 The nearest neighbor search identifies an index set $\mathbb{I}_k(\mathbf{x}_T^{(j)}) \subset \mathbb{I}_C$ of k nearest neighbors for every target
 343 gene count vector $\mathbf{x}_T^{(j)}$. That is, for every context cell with index $i \in \mathbb{I}_k(\mathbf{x}_T^{(j)})$, the chosen measure of
 344 association¹ between the homologous gene counts $\mathbf{h}_C^{(i)}$ and $\mathbf{h}_T^{(j)}$ is lower than for cells outside the set:

$$d(\mathbf{h}_C^{(i)}, \mathbf{h}_T^{(j)}) \leq d(\mathbf{h}_C^{(l)}, \mathbf{h}_T^{(j)}) \text{ for all } l \in \mathbb{I}_C \setminus \mathbb{I}_k(\mathbf{x}_T^{(j)}). \quad (11)$$

345 Common metrics or distance functions can be used as a measure of association d to compare count val-
 346 ues of single-cell data. Some popular choices have been investigated in [30]. We utilize cosine similarity,
 347 measuring the cosine of the angle between $\log 1\text{p}$ -transformed count vectors, as it is fast to calculate even
 348 on datasets containing numerous samples:

$$d(\mathbf{h}_C^{(i)}, \mathbf{h}_T^{(j)}) = 1 - \frac{\langle \log(\mathbf{h}_C^{(i)} + 1), \log(\mathbf{h}_T^{(j)} + 1) \rangle}{\|\log(\mathbf{h}_C^{(i)} + 1)\|_2 \|\log(\mathbf{h}_T^{(j)} + 1)\|_2}. \quad (12)$$

349 The data-level nearest neighbor search can also be used to assign preliminary labels. We count the
 350 multiplicity of cell labels for all context neighbors and assign, as a preliminary label prediction, the most
 351 occurring label,

$$\hat{c}_T^{(j)} = \text{mode} \left[c_C^{(i)} : i \in \mathbb{I}_k(\mathbf{x}_T^{(j)}) \right]. \quad (13)$$

¹Lower values indicate higher association.

352 As the data-level nearest neighbor search is noisy, we additionally assign agreement scores based on the
 353 occurrence of a cell label prediction $\hat{c}_T^{(j)}$.

$$P(\hat{c}_T^{(j)}) = \frac{|\{i : c_C^{(i)} = \hat{c}_T^{(j)} \text{ and } i \in \mathbb{I}_k(\mathbf{x}_T^{(j)})\}|}{k} \quad (14)$$

354 A higher agreement score indicates lower noise, as there is high agreement among cell labels of the
 355 context neighbors. During the following alignment, only target cells exhibiting high agreement scores are
 356 considered for alignment in the intermediate space. For this, we collect all agreement scores for target
 357 cells predicted to have label $\hat{c}_T^{(j)}$ and compute the quantile at level p over this set $\{P(\hat{c}) : \hat{c} = \hat{c}_T^{(j)}\}$. Finally,
 358 we collect the indices of all target cells whose agreement scores of their predicted cell label are higher
 359 than the quantile Q at level p ,

$$\mathbb{J}(p) = \left\{ j : P(\hat{c}_T^{(j)}) > Q\left(p, \{P(\hat{c}) : \hat{c} = \hat{c}_T^{(j)}\}\right) \right\}. \quad (15)$$

360 5.2.3 Aligning the intermediate and latent representations

361 During alignment, the weights of the pretrained encoder part f_{inner}^C are not updated. To guide the model to-
 362 wards leveraging the learned structure, scSpecies aligns intermediate representations with high accuracy
 363 scores

$$\mathbf{t}_T^{(j)} = f_{\text{outer}}^T(\mathbf{x}_T^{(j)}, \mathbf{s}_T^{(j)}), j \in \mathbb{J}(p) \quad (16)$$

364 with a representation of a suitable context neighbor representation

$$\mathbf{t}_C^{(i^*)} = f_{\text{outer}}^C(\mathbf{x}_C^{(i^*)}, \mathbf{s}_C^{(i^*)}), i^* \in \mathbb{I}_k(\mathbf{x}_T^{(j)}). \quad (17)$$

365 This is facilitated by minimizing the squared Euclidean distance.

$$\text{minimize} \left\| \mathbf{t}_T^{(j)} - \mathbf{t}_C^{(i^*)} \right\|_2^2, \text{ if } j \in \mathbb{J}(p). \quad (18)$$

366 The optimal choice $i^* \in \mathbb{I}_k$ for minimization among the k candidates is dynamically determined during the
 367 alignment phase: First, we obtain a set of latent context neighbor variables for the target cells considered
 368 during alignment,

$$\mathbb{I}_k(\mathbf{x}_T^{(j)}) = \left\{ \mathbf{z}_C^{(i)} : i \in \mathbb{I}_k(\mathbf{x}_T^{(j)}) \right\}. \quad (19)$$

369 These latent variables $\mathbf{z}_C^{(i)}$ are then decoded with the batch indicator variable $\mathbf{s}_T^{(j)}$ of their target cell. The
 370 decoder output and target library size $l_T^{(j)}$ parameterize a sampling distribution $P_{\mathbf{x}|\mathbf{z}_C^{(i)}, l_T^{(j)}, \mathbf{s}_T^{(j)}}$, which is

371 used to calculate log density values for every candidate. The cell i^* whose latent representation results
 372 in the highest log density value at $\mathbf{x}_T^{(j)}$ is chosen as optimal neighbor candidate:

$$\mathbf{z}_C^{(i^*)} = \operatorname{argmax}_{\mathbf{z}_C^{(i)} \in \mathbb{L}_k(\mathbf{x}_T^{(j)})} \log \left(p_{\mathbf{X} | \mathbf{z}_C^{(i)}, \mathbf{l}_T^{(j)}, \mathbf{s}_T^{(j)}}(\mathbf{x}_T^{(j)}) \right). \quad (20)$$

373 Using this procedure, it is possible to assign a context neighbor with a fitting cell type if at least one
 374 candidate with this cell type is found in this set. The training criterion for the model scVI^T on the target
 375 dataset for a data point is

$$-\text{ELBO}(\mathbf{x}_T^{(j)}, \mathbf{s}_T^{(j)}, \beta) + \gamma \left\| \mathbf{t}_T^{(j)} - \mathbf{t}_C^{(i^*)} \right\|_2^2 [j \in \mathbb{J}(p)], \quad (21)$$

376 where $[j \in \mathbb{J}(p)]$ is the Iverson Bracket that takes value 1 when an index of a target cell j is in $\mathbb{J}(p)$, and
 377 0 otherwise. This holds true for cells that exhibited a high degree of agreement during the data-level
 378 nearest neighbor search. As minimization in the intermediate space is only incentivized for cells with
 379 these indices, the remaining cells within a mini-batch are grouped around them in a way that minimizes
 380 the nELBO of the scVI model.

381 The scalars $\gamma, \beta \geq 0$ weighing different parts of the loss function, the quantile niveau $p \in [0, 1]$ and number
 382 of nearest neighbors $k \in \mathbb{N}$ are hyperparameters.

383 5.2.4 Transferring cell states and cell types

384 The aligned latent representations $\mathbb{L}_C = \{\mathbf{z}_C^{(i)}\}_{i=1}^{M_C}$ and $\mathbb{L}_T = \{\mathbf{z}_T^{(j)}\}_{j=1}^{M_T}$ can be analyzed for similarities
 385 and differences. For example, their dimensionality can be further reduced into two dimensions using a
 386 dimension reduction algorithm like UMAP [24]. To remove the random influence of the latent sampling
 387 process, we calculate UMAP coordinates using the variational mean parameters μ .

388 We can transfer cell labels or cell states from the context to target species by performing a second
 389 neighbor search on aligned latent representations. A suitable measure of association is the learned log-
 390 density, as it considers the learned manifold of the latent space:

$$d(\mathbf{z}_C^{(i)}, \mathbf{z}_T^{(j)}) = -\log \left(p_{\mathbf{X} | \mathbf{z}_C^{(i)}, \mathbf{l}_T^{(j)}, \mathbf{s}_T^{(j)}}(\mathbf{x}_T^{(j)}) \right) \quad (22)$$

391 We transfer the most common cell type among the top k candidates to the target cell.

392 5.2.5 Comparison of gene profiles

393 To perform a comparison of gene expression profiles between cells of context and target dataset, we tailor
 394 the methods outlined in [31] and [32] to scSpecies. For a latent variable z , we obtain normalized gene

395 expression profiles by decoding it with both decoder networks and averaging over all possible batches \mathbb{S} :

$$\rho_C = \frac{1}{|\mathbb{S}_C|} \sum_{s_C \in \mathbb{S}_C} f_{\text{dec}}^C(\mathbf{z}, s_C), \quad \rho_T = \frac{1}{|\mathbb{S}_T|} \sum_{s_T \in \mathbb{S}_T} f_{\text{dec}}^T(\mathbf{z}, s_T) \quad (23)$$

396 Differences in gene expression profiles can be analyzed for homologous genes, for example, by calculat-
397 ing the log2-fold change (LFC)

$$r_{C,T}^g = \log_2 \left(\frac{\rho_{C,g} + \varepsilon}{\rho_{T,g} + \varepsilon} \right) \quad (24)$$

398 For genes g with low expression levels in both species but still high differences, the offset ε ensures
399 the associated LFC maintains a low order of magnitude. We modify the decoder output layers to avoid
400 artifacts from the softmax function. These artifacts can arise due to highly expressed non-homologous
401 genes or due to different data dimensions. We apply the softmax function to homologous and non-
402 homologous genes separately to obtain

$$\rho_{\text{hom}} = \text{softmax}(\rho_1, \dots, \rho_H), \quad \rho_{\text{nhom}} = \text{softmax}(\rho_{H+1}, \dots, \rho_N), \quad (25)$$

403 where N is the dimensionality of the gene expression vector and H the number of homologous genes.
404 Afterwards, both vectors are scaled so that they sum to one,

$$\rho = \left(\frac{H}{N} \rho_{\text{hom}}^\top, \frac{N-H}{N} \rho_{\text{nhom}}^\top \right)^\top. \quad (26)$$

405 Following [32], for a cell type $C = c_C$ we calculate a mixture distribution of latent states.

$$p_C(\mathbf{z}_C) = \frac{1}{|\mathbb{C}_C(c_C)|} \sum_{\mathbf{x}_C^{(i)} \in \mathbb{C}_C(c_C)} q_{\mathbf{z}|\mathbf{x}_C^{(i)}, s_C^{(i)}}(\mathbf{z}_C) \quad (27)$$

406 The set $\mathbb{C}_C(c_C)$ is the set of cells with label c_C with removed outliers. These outliers are identified
407 by estimating the covariance matrix from variational mean samples μ_C . Cells whose variational mean
408 falls outside the 90%-confidence ellipse described by the covariance estimate are removed. An LFC
409 distribution of homologous genes for cell types present in both datasets can be estimated by sampling
410 latent variables from P_C and computing the corresponding LFC values $r_{C,T}^g$. We calculate the median
411 of the empirical LFC distribution as well as the probability $P(|r_{C,T}^g| > \delta)$ of observing an LFC in gene g
412 higher than level $\delta > 0$.

5.3 Layer-wise relevance propagation

In the following, we briefly describe Layer-wise Relevance Propagation (LRP) [26]. LRP explains the output $f(\mathbf{x})$ of a neural network f by decomposing it into local contributions of input nodes x_i , called relevance scores $R_i(x_i)$ [26]. These relevance scores serve as a measure of each input's influence on the network's output: positive scores ($R_i > 0$) signify a positive influence, whereas negative scores ($R_i < 0$) indicate a negative effect. LRP structurally decomposes the function learned by neural networks into a set of smaller, simpler sub-functions of adjacent layers, while ensuring the conservation of relevance scores across the network. This applies locally, where the sum of the relevance score R_i is conserved across two successive layers of the neural network, and globally between the resulting relevance score for each input node x_i and the output $f(\mathbf{x})$ of the model [26].

Considering a neural network with ReLU activation function, the output a_k of a neuron is given by the input \hat{a}_j of the previous layer and their connected weights w_{jk} of the neurons by

$$a_k = \max\left(0, \sum_j \hat{a}_j w_{jk}\right), \quad (28)$$

including the bias with $\hat{a}_0 = 1$. The relevance scores R_k describe the contribution of each neuron activation \hat{a}_j to a_k . They can be computed by the LRP- γ rule through

$$R_j = \sum_k \frac{\hat{a}_j (w_{jk} + \gamma w_{jk}^+)}{\sum_l \hat{a}_l (w_{lk} + \gamma w_{lk}^+)} R_k. \quad (29)$$

Here, w_{jk}^+ are the positive weights, while γ controls how much these positive contributions are emphasized [33]. LRP methodology aligns with the principles of Deep Taylor Decomposition, which breaks down and redistributes the network's output function $f(\mathbf{x})$ layer by layer through Taylor series expansions. This decomposition allows for the derivation of various LRP rules tailored to the network architecture and the specific function being analyzed [34]. To compute relevance scores for context and target gene expression vectors $\mathbf{x}_C, \mathbf{x}_T$ we propagated the relevance of their latent variational mean parameters $\boldsymbol{\mu}_C, \boldsymbol{\mu}_T$ through the corresponding encoder network. We aggregate relevance scores through averaging over latent dimensions and data points of a cell type. A direct comparison of scores between species is complicated by the influence of non-homologous genes and batch-effects on the relevance scores of homologous genes through the conservation property. Rather, ranked lists of genes by scores can be compared across species.

5.4 Metrics

We evaluated label transfer and clustering performance using four key metrics:

BAS: The balanced accuracy score calculates the proportion of cells correctly labeled in both context and

441 target datasets, averaging over all shared cell types and adjusting for the occurrence of smaller cell
 442 labels by weighing them equally.

443 **ARI:** The adjusted Rand index [35] measures the similarity between predicted and true cell labels, cor-
 444 recting for chance. It considers both correct pairings and misclassifications.

445 **AMI:** The adjusted mutual information [35] quantifies how much information the predicted labels share
 446 with the true labels, adjusting for random label assignments.

447 **DBI:** The Davies-Bouldin index [36] evaluates clustering quality by comparing the compactness of clus-
 448 ters to the separation between them. Lower values indicate better clustering.

449 These metrics collectively assess the accuracy of cell type label transfer and the quality of cell clustering
 450 in the aligned latent space. Details regarding their calculation are found in the documentation of the
 451 package `skikit learn` [37] which we used to calculate these metrics.

452 5.5 Hyperparameters

Model	Layer	In	Architecture	Out
f_{outer}	1	$N + S$	Linear, LN, ReLU, Dropout	300
f_{inner}	1	300	Linear, LN, ReLU, Dropout	200
	2	200	Linear $\rightarrow 2 \cdot 10$ Rep. trick	10
$f_{\text{enc L}}$	1	$N + S$	Linear, LN, ReLU, Dropout	200
	2	200	Linear $\rightarrow 2 \cdot 1$ Rep. trick	1
f_{dec}	1	$10 + S$	Linear, LN, ReLU, Dropout	200
	2	200	Linear, LN, ReLU, Dropout	300
	3	300	Linear, (Softmax, Sigmoid)	$2N$
$\theta_{g,s}$		S	Matrix multiplication	N

Table 1: The network architecture used for all models. N denotes the gene expression data dimension, and S the number of batch effects. Layer functions contain an affine linear transformation, followed by layer normalization (LN), ReLU activation functions which are clipped to the interval $[0, 6]$, and dropout layers with a dropout rate of $p = 0.1$. Latent representations are obtained from the variational mean and scale encoder model output via the reparametrization trick.

453 All models were trained with the same network architecture. Gene expression was modeled using
 454 a zero-inflated negative binomial distribution with constant dispersion for genes within an experimental
 455 batch. We chose a 10-dimensional latent space and a 300-dimensional intermediate space and mapped
 456 to and from these spaces with network architectures listed in Table 1. We trained models for 30 epochs
 457 on datasets with more than 10,000 cells and 60 epochs on datasets with less observed samples. Network
 458 parameters were updated with the ADAM optimizer [38] using standard hyperparameters and a batch size
 459 of $M = 128$.

460 We chose to weigh the KL-Divergence terms with $\beta = 0.1$ at epoch 1, incrementally increasing their influ-
461 ence to $\beta = 1$ over 10 epochs. Similarly, the alignment term started with a weight of $\eta = 10$, which was
462 raised to $\eta = 25$. The number of nearest neighbors was set to $k = 25$ and the quantile cut-off for align-
463 ment was set to $p = 0.8$ across datasets exceeding 10,000 samples. For smaller datasets, we lowered
464 the threshold to $p = 0.6$ to avoid discrimination against scarce cell types. In the latent nearest neighbor
465 search, we pre-computed for each target cell a set of 200 nearest neighbors using the Euclidean distance
466 between the variational mean vectors. Among the 25 cells that resulted in the highest likelihood values,
467 we transferred the most occurring cell label. For differential gene expression analysis, we sampled 10,000
468 times from the plugin estimator and set the offset variable to $\varepsilon = 10^{-6}$.

469 To compute layer-wise relevance scores we retrained the networks with unbounded ReLU activation func-
470 tions and without layer normalization, as it is difficult for LRP to handle normalization layers. To counteract
471 exploding intermediate values caused by high gene expression values, we trained the model on log1p-
472 transformed values. Omitting layer normalization lead to a slight performance drop of around 2.5% across
473 all performance metrics. We calculated relevance scores using the LRP- γ rule with $\gamma = 0.15$.

474 We trained both scArches and scPoli on a scVI base model using the scArches package implementa-
475 tion. These models were trained with the same network architecture as scSpecies. We trained both
476 models on homologous genes, as the scArches publication states that zero-filling only produces reliable
477 results when less than 25% of genes are affected [9][See feature overlap between reference and query].
478 scPoli received training with 10-dimensional batch representations. All other hyperparameters were left
479 at default values.

480 5.6 Pre-processing of the datasets

481 Our model underwent testing on publicly available datasets. (Table 2)

482 The 'Liver Cell Atlas' [17, 39] contains a diverse collection of liver cells from multiple species, including
483 mice (both with and without non-alcoholic fatty liver disease), humans, pigs, monkeys, chickens, and
484 hamsters. We utilized all cells acquired through the scRNA-seq and CITE-seq pipelines.

485 The 'Single-Cell Atlas of Human and Mouse White Adipose Tissue' [18, 40] contains gene expression
486 data from human and murine white fat cells. We selected cell samples obtained via single-nucleus se-
487 quencing.

488 The 'Brain Immune Atlas' profiles immune response to a grade IV glioma. For humans we selected cells
489 obtained via scRNA-seq of newly diagnosed and recurrent glioblastoma. For mice we selected cells from
490 the immune response to transplanted glioblastoma [19, 41].

491 We applied a uniform pre-processing pipeline across all datasets. Initially, the dimension of gene expres-
492 sion vectors was reduced to 4000 most highly variable genes [42]. Then we excluded cells with less than
493 2% nonzero genes or belonging to extremely scarce batch and cell labels with less than 20 samples. To

Dataset	Organism	Shared genes	Cells		Batches	Number of cell types	
			<i>H</i>	<i>M</i>		<i>S</i>	Coarse
Liver	<i>C</i> Mouse	4 000	165 680	34	15 (15)	36 (36)	
	<i>T</i> Mouse NAFLD	2 860	91 787	22	14 (14)	28 (22)	
	<i>T</i> Human	1 808	146 839	30	15 (14)	32 (20)	
	<i>T</i> Human small	1 808	5 000	30	15 (14)	32 (20)	
	<i>T</i> Pig	1 694	21 907	2	9 (8)	unknown	
	<i>T</i> Monkey	1 293	8 483	2	7 (7)	unknown	
	<i>T</i> Chicken	1 197	7 456	2	9 (7)	unknown	
	<i>T</i> Hamster	1 662	5 955	2	11 (9)	unknown	
White fat	<i>C</i> Mouse	4 000	192 470	26	17 (17)	47 (47)	
	<i>T</i> Human	1 937	137 306	24	16 (15)	44 (37)	
Glioblastoma	<i>C</i> Mouse	4 000	46 321	6	14 (14)	23 (23)	
	<i>T</i> Human	1 823	58 560	12	14 (14)	24 (22)	

Table 2: The datasets employed for evaluating scSpecies use mice as context species *C*. The number *H* of homologous genes of context and target dataset are listed in the third column. Furthermore, all datasets are annotated with cell type labels, both at coarse and fine levels. The amount of distinct labels are detailed in the 'Number of cell labels' columns. Additionally, the amount of shared cell labels with the context dataset, are indicated in parentheses.

494 obtain a consistent nomenclature between the datasets some cell labels were renamed. In the liver and
495 glioblastoma datasets, some cells have inconsistent cell type labels. For example, some human liver cells
496 are labeled as neutrophils in the fine and monocytes in the coarse cell label category. We excluded all
497 cells with such a labeling conflict.

498 References

- 499 1. Leonelli, S. & Ankeny, R. A. What makes a model organism? *Endeavour* **37**, 209–212. ISSN: 0160-
500 9327. <https://www.sciencedirect.com/science/article/pii/S0160932713000379> (2013).
- 501 2. Canales, C. P. & Walz, K. in *Cellular and Animal Models in Human Genomics Research* (eds Walz,
502 K. & Young, J. I.) 119–140 (Academic Press, 2019). ISBN: 978-0-12-816573-7. <https://www.sciencedirect.com/science/article/pii/B9780128165737000067>.
- 503 3. F, M., L, M. & RD., C. From mice to humans. *Current diabetes reports vol. 12* (2012).
- 504 4. Haddad, A. F. *et al.* Mouse models of glioblastoma for the evaluation of novel therapeutic strategies.
505 *Neuro-Oncology Advances* **3**, vdab100. ISSN: 2632-2498. eprint: [https://academic.oup.com/noa/](https://academic.oup.com/noa/article-pdf/3/1/vdab100/40080542/vdab100.pdf)
506 [article-pdf/3/1/vdab100/40080542/vdab100.pdf](https://academic.oup.com/noa/article-pdf/3/1/vdab100/40080542/vdab100.pdf). <https://doi.org/10.1093/noajnl/vdab100>
507 (July 2021).
- 508 5. Lau, J. K. C., Zhang, X. & Yu, J. Animal models of non-alcoholic fatty liver disease: current perspec-
509 tives and recent advances. en. *J. Pathol.* **241**, 36–44 (Jan. 2017).
- 510 6. Stripecke, R. *et al.* Innovations, challenges, and minimal information for standardization of human-
511 ized mice. en. *EMBO Mol. Med.* **12**, e8662 (July 2020).

- 513 7. Cao, Z.-J., Wei, L., Lu, S., Yang, D.-C. & Gao, G. Searching large-scale scRNA-seq databases via
514 unbiased cell embedding with Cell BLAST. *Nature Communications* **11**. ISSN: 2041-1723. <http://dx.doi.org/10.1038/s41467-020-17281-7> (July 2020).
515
- 516 8. Hu, J. *et al.* Iterative transfer learning with neural network for clustering and cell type classification
517 in single-cell RNA-seq analysis. *Nature Machine Intelligence* **2**, 607–618. ISSN: 2522-5839. <http://dx.doi.org/10.1038/s42256-020-00233-7> (Oct. 2020).
518
- 519 9. Lotfollahi, M. *et al.* Mapping single-cell data to reference atlases by transfer learning. *Nature Biotech-*
520 *nology*, 1–10 (2021).
- 521 10. De Donno, C. *et al.* Population-level integration of single-cell datasets enables multi-scale analysis
522 across samples. *Nature Methods* **20**, 1683–1692. ISSN: 1548-7105. <https://doi.org/10.1038/s41592-023-02035-2> (Nov. 2023).
523
- 524 11. Lotfollahi, M. *et al.* Biologically informed deep learning to query gene programs in single-cell atlases.
525 *Nature Cell Biology* **25**, 337–350. <https://doi.org/10.1038/s41556-022-01072-x> (2023).
- 526 12. Michielsen, L. *et al.* Single-cell reference mapping to construct and extend cell-type hierarchies.
527 *NAR Genomics and Bioinformatics* **5**, lqad070. ISSN: 2631-9268. eprint: <https://academic.oup.com/nargab/article-pdf/5/3/lqad070/51052048/lqad070.pdf>. <https://doi.org/10.1093/nargab/lqad070> (July 2023).
528
529
- 530 13. Breschi, A., Gingeras, T. R. & Guigó, R. Comparative transcriptomics in human and mouse. *Nature*
531 *Reviews Genetics* **18**, 425–440. ISSN: 1471-0064. <https://doi.org/10.1038/nrg.2017.19> (July
532 2017).
- 533 14. Rosen, Y. *et al.* Toward universal cell embeddings: integrating single-cell RNA-seq datasets across
534 species with SATURN. *Nature Methods* **21**, 1492–1500. [https://doi.org/10.1038/s41592-024-](https://doi.org/10.1038/s41592-024-02191-z)
535 [02191-z](https://doi.org/10.1038/s41592-024-02191-z) (Aug. 1, 2024).
- 536 15. Biharie, K., Michielsen, L., Reinders, M. J. T. & Mahfouz, A. Cell type matching across species using
537 protein embeddings and transfer learning. *Bioinformatics* **39**, i404–i412. ISSN: 1367-4811. eprint:
538 [https://academic.oup.com/bioinformatics/article-pdf/39/Supplement_1/i404/50741455/](https://academic.oup.com/bioinformatics/article-pdf/39/Supplement_1/i404/50741455/btad248_supplementary_data.pdf)
539 [btad248_supplementary_data.pdf](https://academic.oup.com/bioinformatics/article-pdf/39/Supplement_1/i404/50741455/btad248_supplementary_data.pdf). <https://doi.org/10.1093/bioinformatics/btad248>
540 (June 2023).
- 541 16. Sohn, K., Yan, X. & Lee, H. *Learning Structured Output Representation Using Deep Conditional*
542 *Generative Models* in *Proceedings of the 28th International Conference on Neural Information Pro-*
543 *cessing Systems - Volume 2* (MIT Press, Montreal, Canada, 2015), 3483–3491.
- 544 17. Guilliams, M. *et al.* Spatial proteogenomics reveals distinct and evolutionarily conserved hepatic
545 macrophage niches. *Cell* **185**, 379–396 (Jan. 2022).

- 546 18. Emont, M. P. *et al.* A single-cell atlas of human and mouse white adipose tissue. *Nature* **603**, 926–
547 933. ISSN: 1476-4687. <https://doi.org/10.1038/s41586-022-04518-2>, (Mar. 2022).
- 548 19. Pombo Antunes, A. R. *et al.* Single-cell profiling of myeloid cells in glioblastoma across species and
549 disease stage reveals macrophage competition and specialization. *Nature Neuroscience* **24**, 595–
550 610. <https://doi.org/10.1038/s41593-020-00789-y> (Apr. 1, 2021).
- 551 20. Lopez, R., Regier, J., Cole, M., Jordan, M. I. & Yosef, N. Deep Generative Modeling for Single-cell
552 Transcriptomics. *Nature methods* **15**, 1053–1058. [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:53643161)
553 53643161 (2018).
- 554 21. Fernando, B., Fromont, E. & Tuytelaars, T. Mining Mid-level Features for Image Classification. *Inter-*
555 *national Journal of Computer Vision* **108**, 186–203. ISSN: 1573-1405. [https://doi.org/10.1007/](https://doi.org/10.1007/s11263-014-0700-1)
556 s11263-014-0700-1 (July 2014).
- 557 22. Boureau, Y.-L., Bach, F., LeCun, Y. & Ponce, J. *Learning mid-level features for recognition in 2010*
558 *IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (2010), 2559–
559 2566.
- 560 23. Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. *How transferable are features in deep neural net-*
561 *works?* in *Proceedings of the 27th International Conference on Neural Information Processing Sys-*
562 *tems - Volume 2* (MIT Press, Montreal, Canada, 2014), 3320–3328.
- 563 24. McInnes, L., Healy, J. & Melville, J. *UMAP: Uniform Manifold Approximation and Projection for Di-*
564 *mension Reduction* 2020. arXiv: 1802.03426 [stat.ML].
- 565 25. Jiang, C. *et al.* Comparative transcriptomics analyses in livers of mice, humans, and humanized
566 mice define human-specific gene networks. en. *Cells* **9**, 2566 (Nov. 2020).
- 567 26. Bach, S. *et al.* On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Rele-
568 vance Propagation. *PLOS ONE*, 46 (2015).
- 569 27. Keyl, P. *et al.* Single-cell gene regulatory network prediction by explainable AI. en. *Nucleic Acids*
570 *Res.* **51**, e20 (Feb. 2023).
- 571 28. Ma, X.-Y. *et al.* Malat1 as an evolutionarily conserved lncRNA, plays a positive role in regulating
572 proliferation and maintaining undifferentiated status of early-stage hematopoietic cells. en. *BMC*
573 *Genomics* **16**, 676 (Sept. 2015).
- 574 29. Kingma, D. P. & Welling, M. *Auto-Encoding Variational Bayes* 2022. arXiv: 1312.6114 [stat.ML].
- 575 30. Skinnider, M. A., Squair, J. W. & Foster, L. J. Evaluating measures of association for single-cell
576 transcriptomics. *Nature methods* **16**, 381–386. ISSN: 1548-7091. [https://doi.org/10.1038/](https://doi.org/10.1038/s41592-019-0372-4)
577 s41592-019-0372-4 (May 2019).

- 578 31. Boyeau, P. *et al.* Deep Generative Models for Detecting Differential Expression in Single Cells.
579 *bioRxiv*. eprint: <https://www.biorxiv.org/content/early/2019/10/04/794289.full.pdf>.
580 <https://www.biorxiv.org/content/early/2019/10/04/794289> (2019).
- 581 32. Boyeau, P. *et al.* An empirical Bayes method for differential expression analysis of single cells with
582 deep generative models. *Proceedings of the National Academy of Sciences* **120**, e2209124120.
583 eprint: <https://www.pnas.org/doi/pdf/10.1073/pnas.2209124120>. <https://www.pnas.org/doi/abs/10.1073/pnas.2209124120> (2023).
584
- 585 33. Montavon, G., Binder, A., Lapuschkin, S., Samek, W. & Müller, K.-R. in *Explainable AI: Interpreting,*
586 *Explaining and Visualizing Deep Learning* (eds Samek, W., Montavon, G., Vedaldi, A., Hansen, L. K.
587 & Müller, K.-R.) Series Title: Lecture Notes in Computer Science, 193–209 (Springer International
588 Publishing, Cham, 2019). http://link.springer.com/10.1007/978-3-030-28954-6_10 (2022).
- 589 34. Montavon, G., Bach, S., Binder, A., Samek, W. & Müller, K.-R. Explaining NonLinear Classifica-
590 tion Decisions with Deep Taylor Decomposition. *Pattern Recognition* **65**, 211–222. ISSN: 00313203.
591 arXiv: 1512.02479[cs,stat]. <http://arxiv.org/abs/1512.02479> (2022) (May 2017).
- 592 35. Vinh, N. X., Epps, J. & Bailey, J. Information Theoretic Measures for Clusterings Comparison: Vari-
593 ants, Properties, Normalization and Correction for Chance. *J. Mach. Learn. Res.* **11**, 2837–2854.
594 ISSN: 1532-4435 (Dec. 2010).
- 595 36. Davies, D. L. & Bouldin, D. W. A Cluster Separation Measure. *IEEE Transactions on Pattern Analysis*
596 *and Machine Intelligence* **PAMI-1**, 224–227 (1979).
- 597 37. *Userguide to scikit learn* Accessed: 2024-10-01, Balanced accuracy score: Section 3.4.2.4. Adjusted
598 Rand index: Section 2.3.11.1. Adjusted mutual information: Section 2.3.11.2. Davies-Bouldin index:
599 Section 2.3.11.7. https://scikit-learn.org/stable/modules/model_evaluation.html.
- 600 38. Kingma, D. P. & Ba, J. *Adam: A Method for Stochastic Optimization* 2017. arXiv: 1412.6980 [cs.LG].
- 601 39. *Brain Immune Atlas* Accessed: 2023-06-20. <https://www.livercellatlas.org/>.
- 602 40. *Single-Cell Atlas of Human and Mouse White Adipose Tissue* Accessed: 2024-02-15. https://singlecell.broadinstitute.org/single_cell/study/SCP1376.
- 603
- 604 41. *Brain Immune Atlas* Accessed: 2024-03-02. <https://www.brainimmuneatlas.org/>.
- 605 42. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell
606 gene expression data. *Nature Biotechnology* **33**, 495–502. ISSN: 1546-1696. <https://doi.org/10.1038/nbt.3192> (May 2015).
607

608 **6 Declarations**

609 **6.1 Ethics approval and consent to participate**

610 Not applicable.

611 **6.2 Consent for publication**

612 Not applicable.

613 **6.3 Availability of data and materials**

614 The datasets can be accessed via the URLs [39–41].

615 Our model is implemented in Python 3.11.5 with PyTorch 2.1. The preprocessing scripts to obtain the
616 datasets and the code to reproduce our results can be accessed at [https://github.com/cschaech/](https://github.com/cschaech/scSpecies)
617 `scSpecies`. We recommend to use a device equipped with an NVIDIA GPU.

618 **6.4 Competing interests**

619 The authors declare that they have no competing interests.

620 **6.5 Funding**

621 Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID
622 499552394 – SFB 1597 Small Data.

623 **6.6 Authors' affiliation**

624 Institute of Medical Biometry and Statistics, Faculty of Medicine and Medical Center - University of
625 Freiburg, Germany: Clemens Schächter, Martin Treppner, Maren Hackenberg

626 Neurorobotics Lab, Dept. of Computer Science – University of Freiburg, Germany: Hanne Raum, Joschka
627 Bödecker

628 BrainLinks-BrainTools CRIION - Collaborative Research Institute Intelligent Oncology: Joschka Bödecker
629 Freiburg Center for Data Analysis and Modelling – University of Freiburg, Germany: Harald Binder

630 CIBSS, Centre for Integrative Biological Signalling Studies – University of Freiburg, Germany: Harald
631 Binder

632 **6.7 Authors' contributions**

633 H.B. conceived and coordinated the project. H.B., C.S., and M.T. jointly developed the approach for
634 aligning network architectures across species. C.S. implemented the corresponding code. H.R. and J.B.
635 designed the methodology for extending the analysis from the latent space to the data level, with H.R.
636 handling the implementation. C.S., H.B., M.H., and H.R. contributed to the writing of the manuscript. All
637 authors reviewed and approved the final version of the manuscript.

638 **6.8 Correspondence**

639 Correspondence to Clemens Schächter or Harald Binder.

640 **6.9 Acknowledgement**

641 Not applicable

642 **7 Extended Data**

Model	scArches		scPoli		kNN classifier					
	-		-		$k = 1$		$k = 25$		$k = 250$	
Neighbors										
Cell labels	coarse	fine	coarse	fine	coarse	fine	coarse	fine	coarse	fine
Balanced label transfer accuracy score in % (BAS)										
Liver - human	64.05	48.05	80.74	55.72	80.72	59.46	79.70	62.04	75.16	57.25
Liver - mouse	97.62	78.50	98.67	81.44	97.69	79.40	98.03	80.03	97.45	76.08
White fat	65.79	37.50	65.45	37.41	74.37	40.20	73.80	41.17	67.64	37.86
Glioblastoma	51.96	46.60	80.92	59.94	75.59	54.47	76.37	56.65	71.70	54.51
Adjusted Rand index (ARI)										
Liver - human	0.725	0.248	0.841	0.263	0.740	0.194	0.824	0.253	0.859	0.290
Liver - mouse	0.983	0.837	0.984	0.825	0.983	0.822	0.985	0.839	0.982	0.844
White fat	0.773	0.414	0.846	0.443	0.868	0.371	0.884	0.438	0.877	0.469
Glioblastoma	0.458	0.401	0.583	0.581	0.481	0.384	0.537	0.455	0.525	0.470
Adjusted mutual information (AMI)										
Liver - human	0.685	0.516	0.794	0.538	0.711	0.487	0.781	0.554	0.809	0.575
Liver - mouse	0.976	0.871	0.983	0.869	0.977	0.860	0.981	0.875	0.977	0.870
White fat	0.768	0.607	0.831	0.657	0.839	0.599	0.861	0.654	0.848	0.659
Glioblastoma	0.576	0.500	0.656	0.598	0.610	0.507	0.679	0.568	0.672	0.568
scSpecies										
	lat. alignment		intermediate alignment							
Neighbors	$k = 25$		$k = 0$		$k = 1$		$k = 25$		$k = 250$	
Cell labels	coarse	fine	coarse	fine	coarse	fine	coarse	fine	coarse	fine
Balanced label transfer accuracy score in % (BAS)										
Liver - human	90.35	71.12	5.01	2.81	86.35	66.74	92.08	73.29	91.54	71.62
Liver - small	86.57	65.67	7.91	4.52	79.45	59.59	87.76	67.78	81.19	62.66
Liver - mouse	97.56	80.40	5.36	1.83	97.99	81.06	98.11	81.24	97.82	79.51
White fat	79.31	48.81	5.79	2.27	78.14	47.02	82.02	49.15	83.17	48.42
Glioblastoma	88.41	67.54	9.61	6.26	84.69	63.87	88.88	68.87	84.07	64.90
Adjusted Rand index (ARI)										
Liver - human	0.865	0.456	0.204	0.163	0.872	0.406	0.888	0.509	0.887	0.593
Liver - small	0.841	0.451	0.237	0.181	0.747	0.275	0.863	0.481	0.849	0.545
Liver - mouse	0.975	0.832	0.182	0.192	0.985	0.834	0.987	0.837	0.984	0.834
White fat	0.944	0.519	0.142	0.137	0.880	0.487	0.959	0.528	0.963	0.540
Glioblastoma	0.717	0.648	0.144	0.216	0.633	0.551	0.753	0.684	0.734	0.666
Adjusted mutual information (AMI)										
Liver - human	0.824	0.703	0.351	0.408	0.827	0.673	0.855	0.731	0.864	0.760
Liver - small	0.805	0.676	0.334	0.354	0.697	0.540	0.825	0.696	0.830	0.727
Liver - mouse	0.971	0.870	0.380	0.455	0.980	0.875	0.981	0.878	0.978	0.876
White fat	0.912	0.711	0.268	0.352	0.867	0.690	0.929	0.725	0.934	0.734
Glioblastoma	0.782	0.698	0.246	0.401	0.745	0.628	0.799	0.683	0.783	0.675

Table 3: Comparison of model performance on four different datasets. The results are averaged over five random seeds and the best results highlighted by bold font. The results for each dataset are listed for the coarse - fine cell label categories. The upper table contains the results obtained by scArches and scPoli. The kNN columns refer to the results of a data-level k nearest neighbor classifier trained on shared homologous genes. The results from scSpecies are listed in the bottom table. The first column corresponds to the results of a scSpecies model where latent representations instead of the intermediate representations are aligned. The column with zero neighbors corresponds to completely omitting the nearest neighbor integration within the model. The column with one neighbor corresponds to omitting learning a suitable neighbor candidate, as the choice is fixed.

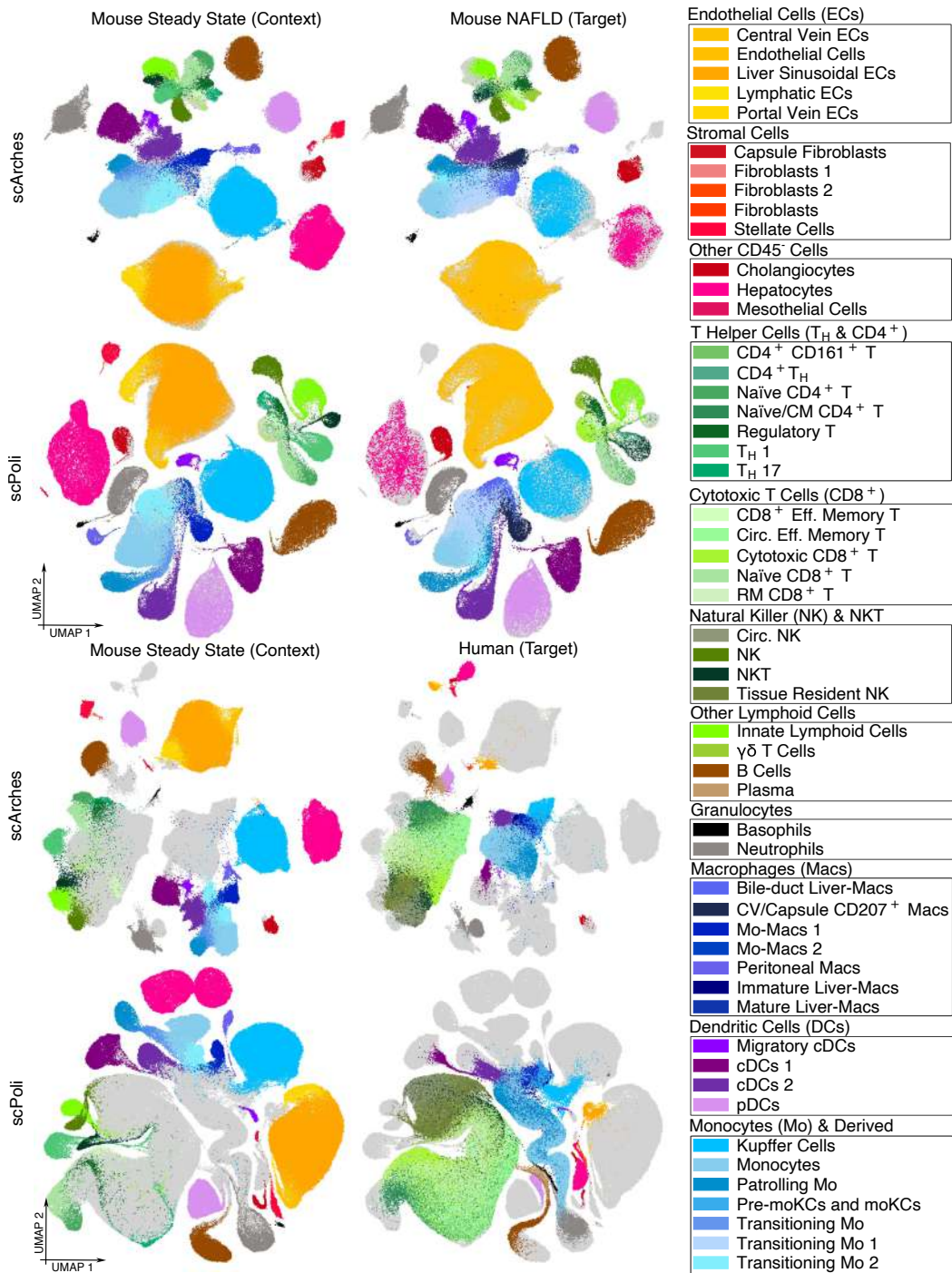


Figure 4: Alignment performance of the architecture surgery-based approaches scArches and scPoli. The four left-hand plots were generated by aligning two mouse liver cell datasets. One dataset contains cell samples from healthy organisms, while the other contains cells from mice with non-alcoholic fatty liver disease. Despite the difference in disease conditions the latent representations are well aligned. The four plots on the right side were obtained by aligning human liver cells with those of healthy mice. Here, both approaches encounter difficulties with cross-species alignment.

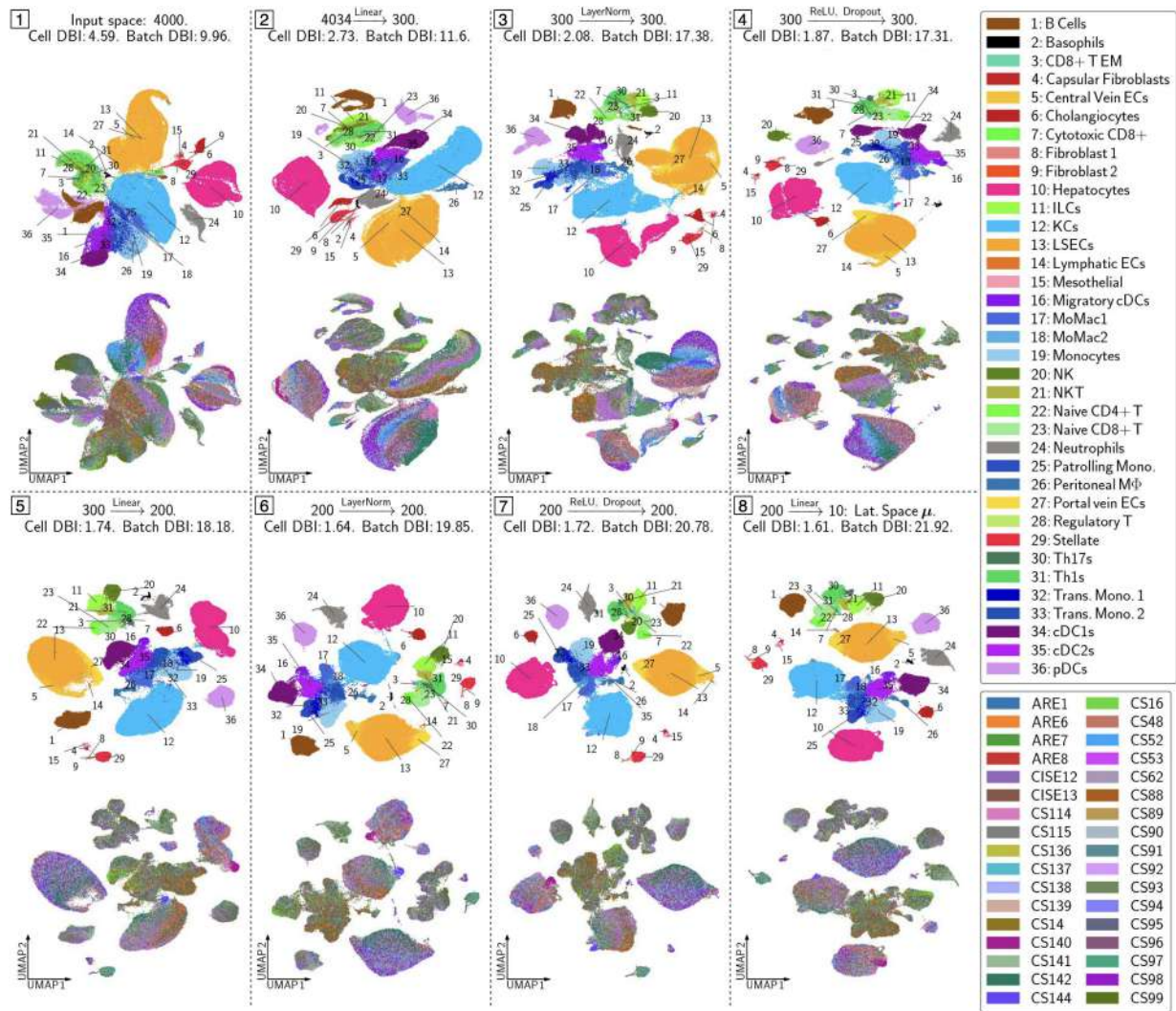


Figure 5: Intermediate spaces of a scVI model applied to the mouse liver context dataset. It details the layer transformations from data space to latent space. Subplot 1 represents the UMAP coordinates of the original dataset, while subplot 8 shows the variational mean vectors in the latent space. Subplots 2–7 depict the UMAP coordinates of the intermediate dataset representation obtained by applying the corresponding layer transformation. Each subplot presents two scatter plots: the upper one showing clusters based on cell labels and the lower one depicting experimental batches. Additionally, the Davies-Bouldin index is used to assess the clustering quality for each subplot.

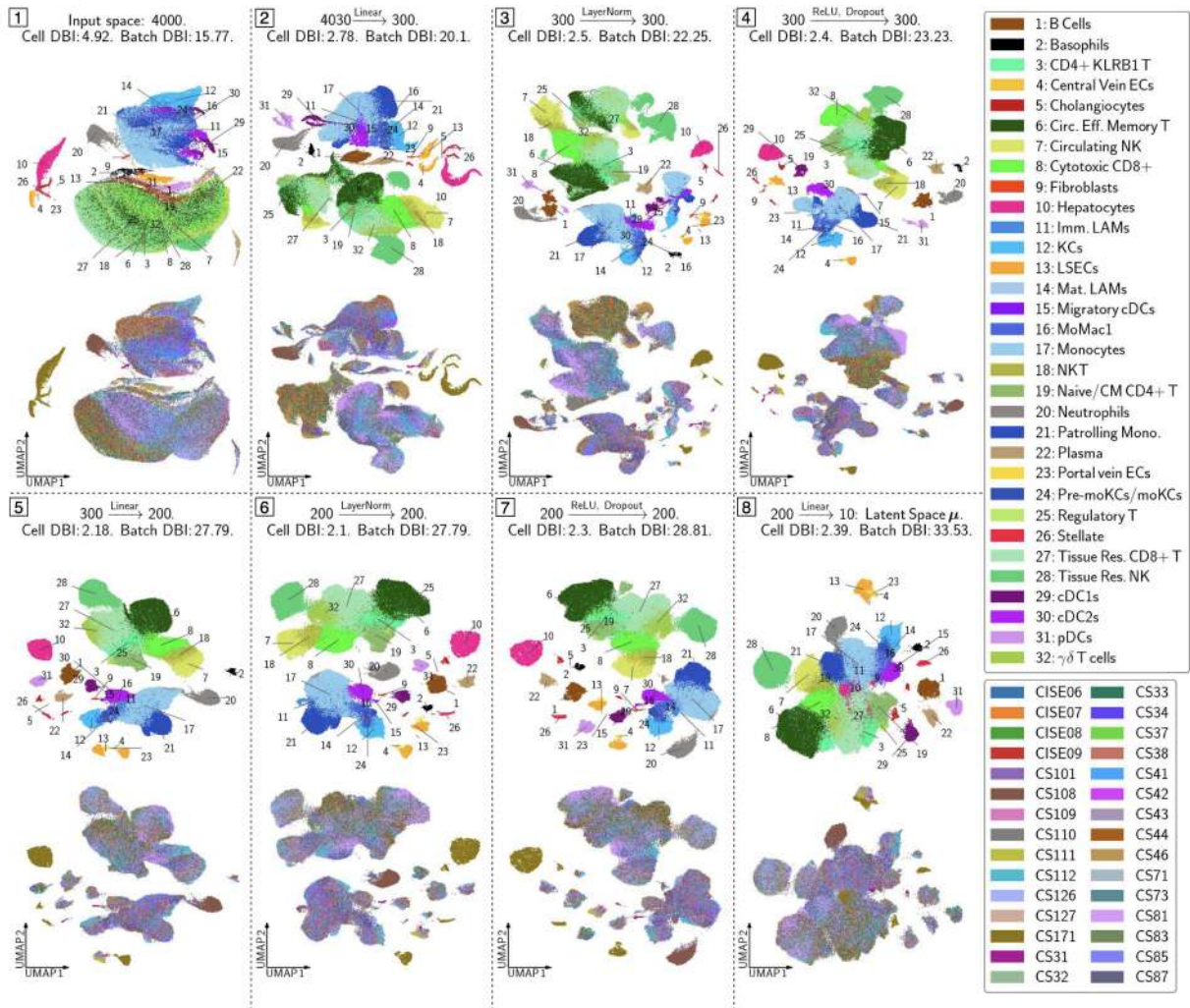


Figure 6: Intermediate spaces of a scVI model applied to the unaligned human liver target dataset. For an explanation of the subplots, see Figure 5.

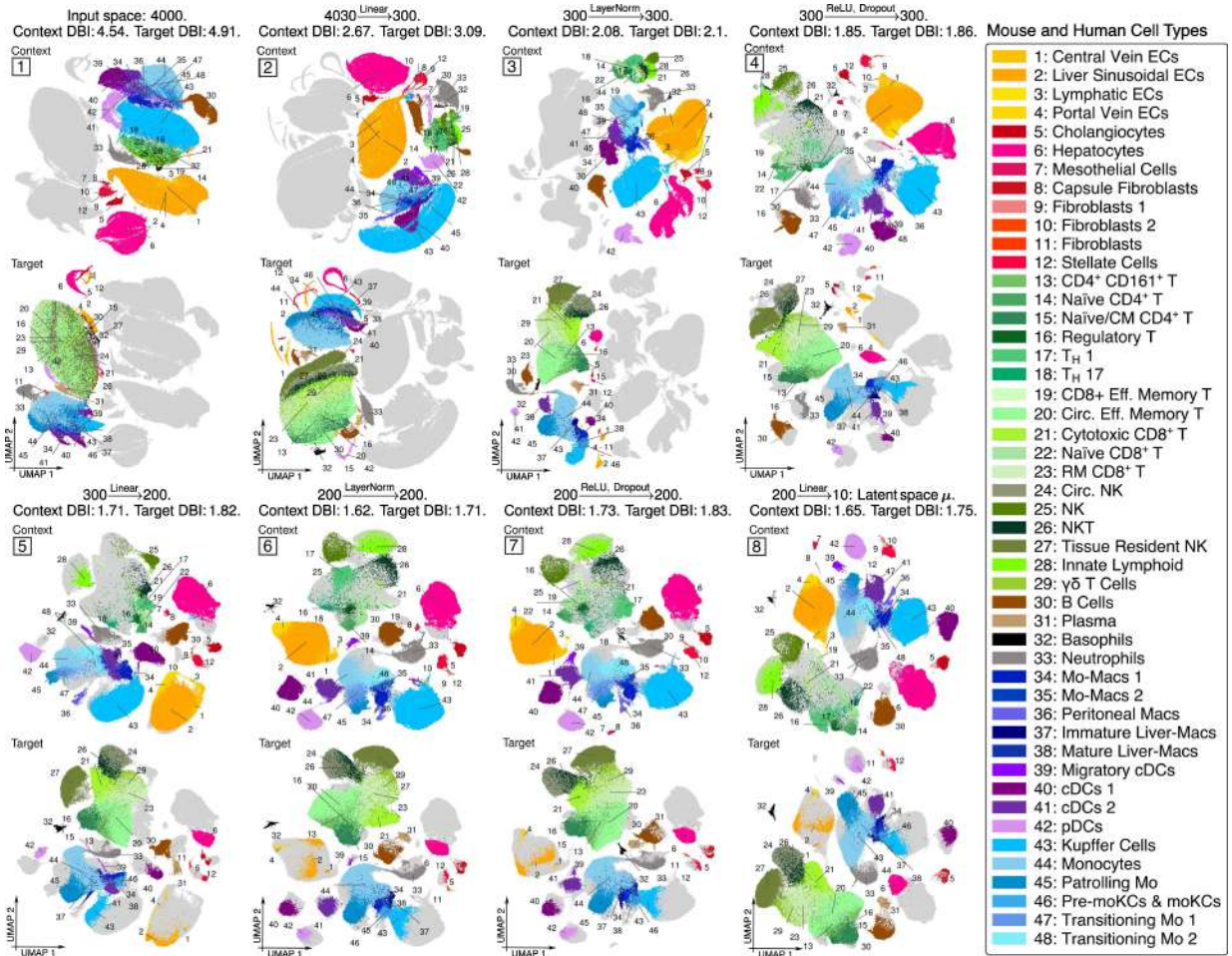


Figure 7: Intermediate spaces of a scSpecies model applied to the mouse-human liver dataset pair. Each subplot presents two scatter plots: the upper one showing context cell label clusters and the lower one depicting the human target cell clusters. Additionally, the Davies-Bouldin index is used to assess clustering quality for each subplot. Alignment of the two datasets is encouraged in subplot 4.

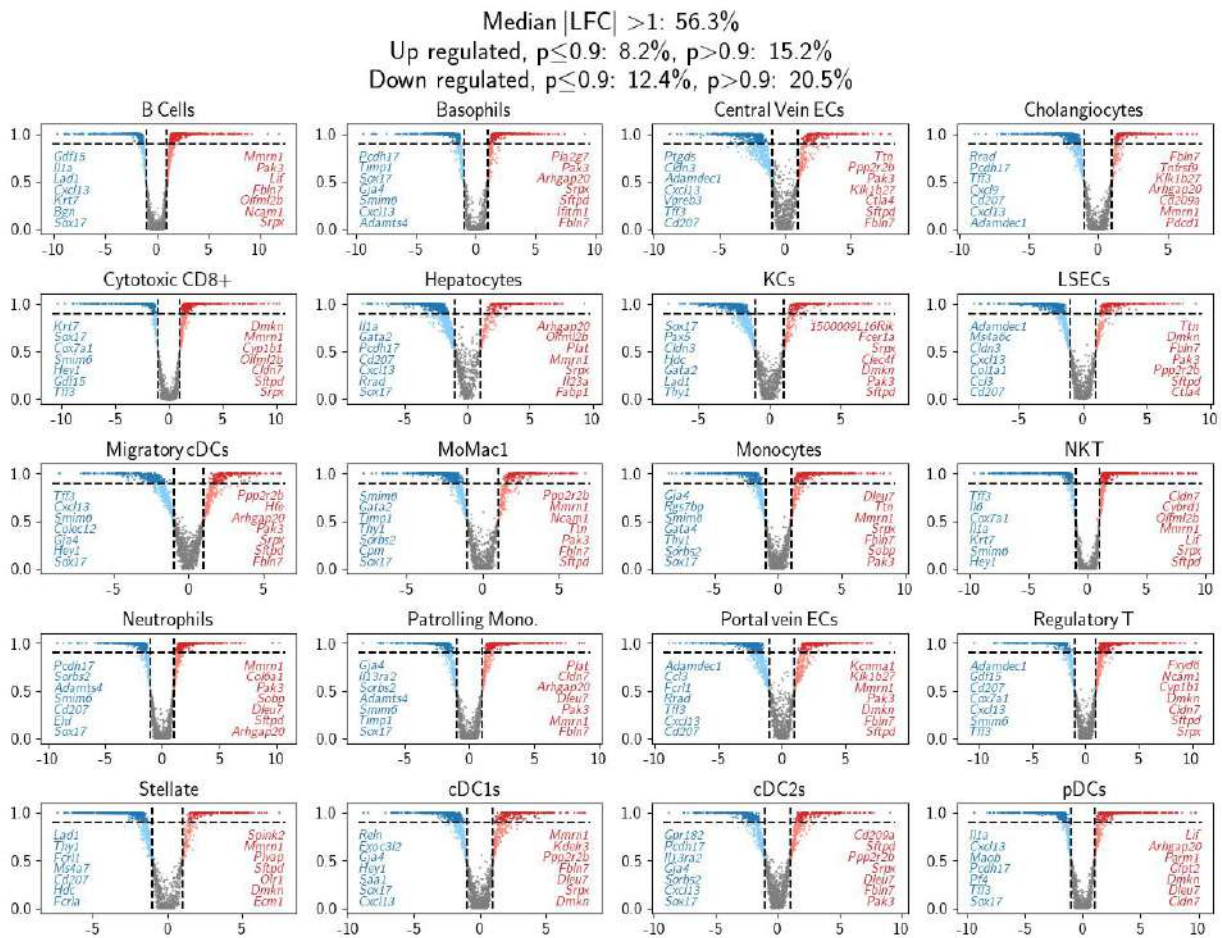


Figure 8: A comparative analysis of gene expression profiles between humans and mice using scSpecies. We computed the median of the empirical log₂ fold change distribution, displayed along the x-axis. The y-axis illustrates the likelihood of a gene being differentially expressed in mice versus humans with an LFC exceeding one. The compared cells are decoded from a randomly selected latent value within a latent cell type distribution. The figure highlights the top seven genes in mice that are significantly up-regulated (indicated in red) and the top seven that are notably down-regulated (blue) in comparison to their human equivalents.

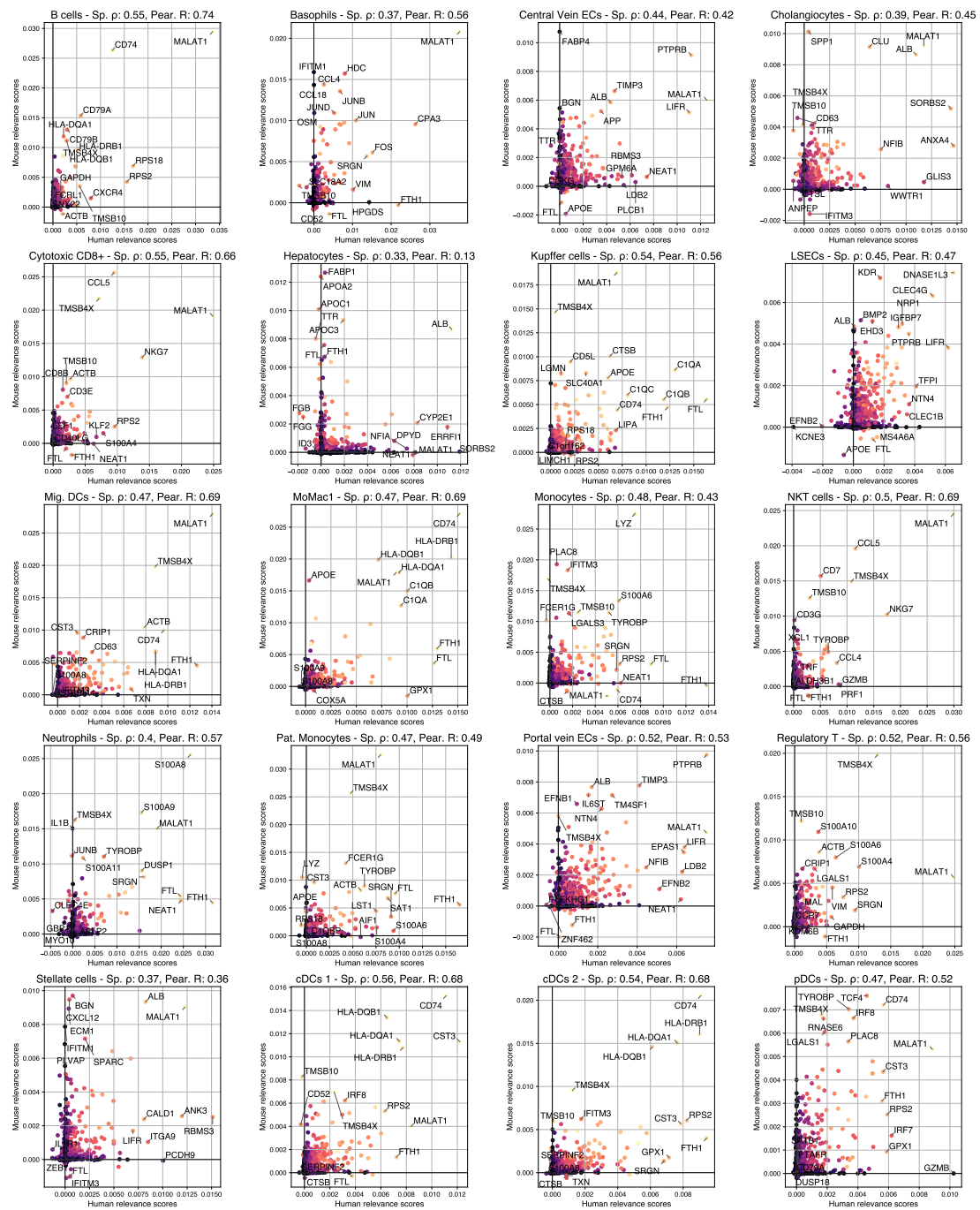


Figure 9: Plots of human and mouse gene LRP scores against each other. Each dot represents a homologous gene. For every cell, Spearman's ρ and Person's R between human and mice LRP values are given in the axis label. Coloring corresponds to combined products of human and mice gene expression, with values of 0 are colored in dark tones and high values in bright colors.

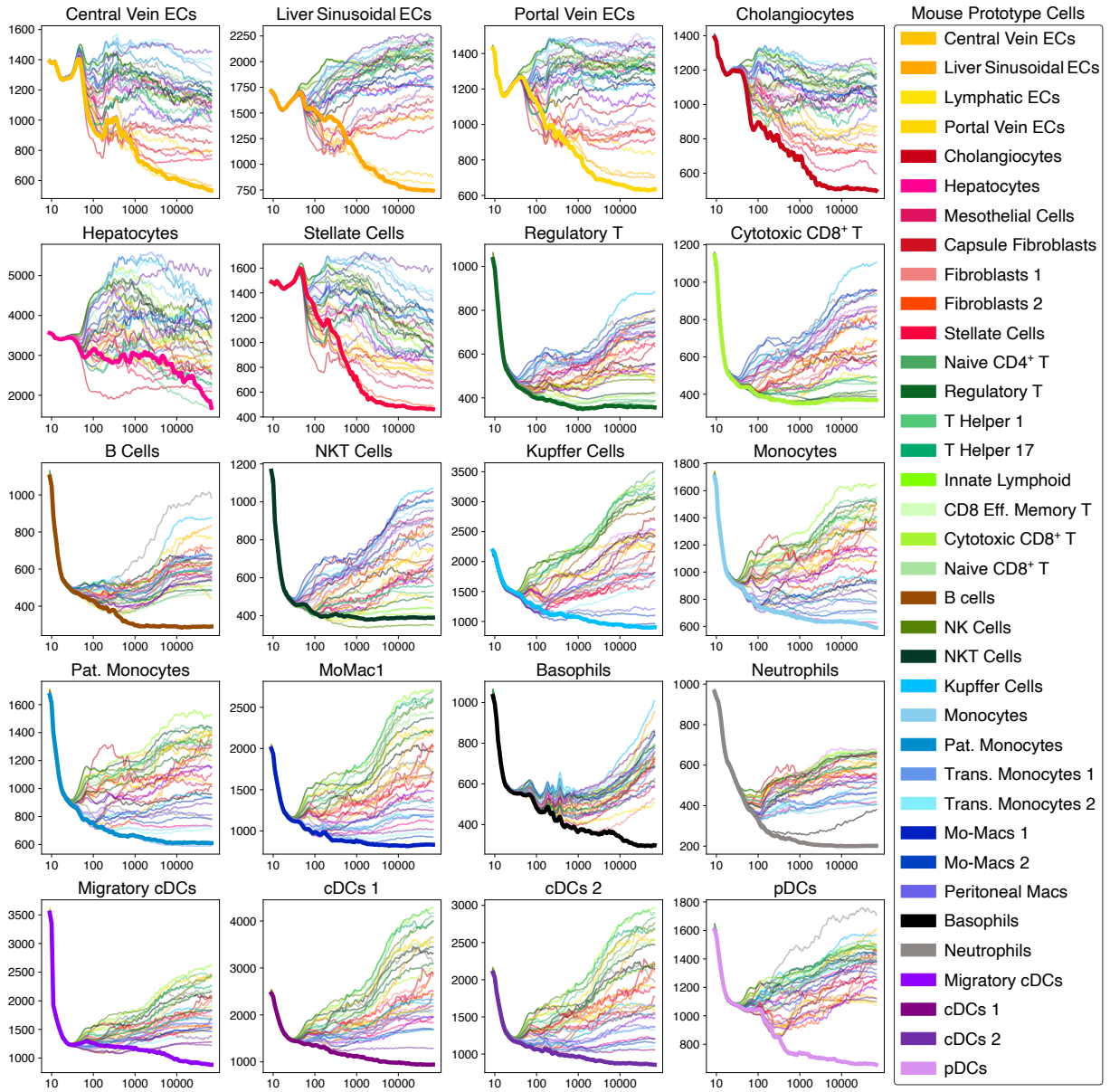


Figure 10: Illustration of the alignment process of scSpecies with $k = 25$ neighbors. On the y-axis, we plot the negative log-density values derived from reconstructing human liver cell prototypes using their candidate set of mouse latent variables. The x-axis shows a log-scale trajectory of these values, averaged over the last $\lceil \min(10, 0.05 \times \text{steps}) \rceil$ iterations.