# Identifying Reproducible Transcription Regulator Coexpression Patterns with Single Cell Transcriptomics

*Alexander Morin[1,2,3], C. Pan Chu[1,2,3], Paul Pavlidis[1,2,*]*

1. Michael Smith Laboratories, University of British Columbia, Vancouver, BC, Canada

2. Department of Psychiatry, University of British Columbia, Vancouver, BC, Canada

3. Graduate Program in Bioinformatics, University of British Columbia, Vancouver, BC, Canada

* Corresponding author

Paul Pavlidis

177 Michael Smith Laboratories

2185 East Mall

University of British Columbia

Vancouver BC V6T1Z4

Canada

604 827 4157

paul@msl.ubc.ca

# Abstract

20 

21  The proliferation of single cell transcriptomics has potentiated our ability to unveil
22  patterns that reflect dynamic cellular processes such as the regulation of gene
23  transcription. In this study, we leverage a broad collection of single cell RNA-seq data to
24  identify the gene partners whose expression is most coordinated with each human and
25  mouse transcription regulator (TR). We assembled 120 human and 103 mouse scRNA-
26  seq datasets from the literature (>28 million cells), constructing a single cell
27  coexpression network for each. We aimed to understand the consistency of TR
28  coexpression profiles across a broad sampling of biological contexts, rather than
29  examine the preservation of context-specific signals. Our workflow therefore explicitly
30  prioritizes the patterns that are most reproducible across cell types. Towards this goal,
31  we characterize the similarity of each TR's coexpression within and across species. We
32  create single cell coexpression rankings for each TR, demonstrating that this
33  aggregated information recovers literature curated targets on par with ChIP-seq data.
34  We then combine the coexpression and ChIP-seq information to identify candidate
35  regulatory interactions supported across methods and species. Finally, we highlight
36  interactions for the important neural TR ASCL1 to demonstrate how our compiled
37  information can be adopted for community use.

# Author Summary

38 

39  A common way to analyze gene expression (transcriptomics) data is to correlate gene
40  transcript levels across samples for every pair of genes (coexpression). Coordinated
41  expression between genes may imply a shared biological function, though this warrants
42  cautious interpretation given assumptions about cellular processes inferred from RNA
43  abundances alone. Still, coexpression inference is often used to nominate genes whose
44  expression may be controlled by transcription regulators (TRs). The rapid generation of
45  diverse single cell transcriptomics data has unlocked our ability to discover
46  coexpression patterns across individual cells — though these signals are often noisy.
47  Reproducible patterns across studies can help distinguish meaningful biological
48  relationships from spurious correlations. We used this study to analyze a broad
49  collection of single cell data spanning numerous tissues in human and mouse to infer
50  global TR coexpression patterns. We aimed to learn which interactions were generally
51  observable, to better potentiate future examinations of reproducible coexpression in
52  specific contexts. We evaluate the predictive performance of these global single cell
53  coexpression rankings using independent gene regulation evidence, and highlight TR-
54  gene pairs that are supported across data modalities as well as species. By
55  disseminating these rankings, we hope that other researchers can extract insight for
56  their own TRs of interest.

# Introduction

The widespread adoption of single cell genomic methodologies, particularly single cell/nucleus RNA sequencing (herein, scRNA-seq), has significantly advanced our ability to characterize dynamic cellular processes. The scale with which scRNA-seq data has been generated has created an unprecedented opportunity to understand the reproducibility of these cellular patterns. This is important because, despite its power, scRNA-seq results in sparse gene transcript counts due to both biological and technical factors (Crow et al., 2016; Heumos et al., 2023).

Gene regulation is a field that stands to greatly benefit from the single cell era. A primary objective is to map the temporal and context-specific interactions between transcription regulators (TRs) and their target genes. However, understanding the sets of genes regulated by each TR — regardless of context — remains a challenge. Despite the availability of genetic tools, linking TRs to direct gene targets is hindered by multiple factors. These include the cost and difficulty of collecting experimental data implicating direct regulation, such as TR binding information from chromatin immunoprecipitation sequencing (ChIP-seq), and the inherent complexity of the underlying biology (Lambert et al., 2018, Rothenberg 2019).

Gene coexpression is a traditional and widely adopted approach for predicting TR-target relationships. This analysis is often cast as generating a predicted gene regulatory network, where the strength of covariation between gene transcript levels serves as edge weights (Sonawane et al., 2019). The fundamental assumption is that if a TR protein influences a gene's transcription, the TR gene itself must also be expressed. However, this assumption may be compromised when the dynamic expression of TRs and their targets are uncoupled. Further, this covariation does not implicate a causative directionality (i.e., regulatory influence) between gene pairs. Despite these limitations, coexpression analysis has been extensively applied as a cost-effective and genome-wide strategy to investigate gene regulation and is commonly integrated with other data modalities (Aibar et al., 2017; Bravo González-Blas et al., 2023).

The emergence of scRNA-seq has made it possible to study coexpression at a finer level of granularity than afforded by bulk tissue, mitigating cell type compositional effects that impact bulk tissue interpretation (McCall et al., 2016; Farahbod and Pavlidis, 2019; Farahbod and Pavlidis, 2020; Zhang et al., 2021). However, cautious interpretation is still warranted due to the sparsity of scRNA-seq data. Correspondingly, the benefits of a meta-analytic framework (Lee et al., 2004; Mistry et al., 2013; Ballouz et al., 2015) have been extended to single cell coexpression to tasks such as gene function prediction (Crow et al., 2016; Crow and Gillis, 2018) and understanding reproducible patterns in the brain (Harris et al. 2021; Suresh et al. 2023; Werner and Gillis 2023). Importantly, these studies typically focused on the preservation of the global coexpression network structure, rather than any specific gene profile.

We drew inspiration from these works and our experience in aggregating ChIP-seq and TR perturbation studies to identify reproducible TR-target interactions (Morin et al., 2023). This stemmed from the recognition that the evidence from various lines of gene

99  regulation methods often do not intersect, necessitating comprehensive data
100 compilation (Hu et al., 2007; Gitter et al., 2009; Cusanovich et al., 2014; Garcia-Alonso
101 et al., 2019; Kang et al., 2020). In this study, we adopt a "TR-centric" approach towards
102 aggregating single cell coexpression networks, with the primary goal of learning
103 reproducible TR interactions. Specifically, our focus was to assemble a diverse range of
104 scRNA-seq data to better understand the coexpression range of all measurable TRs in
105 mouse and human. Our key aim was to prioritize the genes that are most frequently
106 coexpressed with each TR, hypothesizing that this prioritization can facilitate the
107 identification of direct TR-target interactions. We further reasoned that this information
108 would help establish expectations for more focused data aggregations.

# Methods

110 All analyses were performed in the R statistical computing environment (R version 4.2.1
111 http://www.r-project.org). The associated scripts can be found at
112 https://github.com/PavlidisLab/TR_singlecell.

### Genomic tables

114 Gene annotations were based on NCBI RefSeq Select (mm10 and hg38)
115 (https://www.ncbi.nlm.nih.gov/refseq/refseq_select/). High-confidence one-to-one
116 orthologous genes were accessed via the DIOPT resource (V9; Hu et al. 2011; Hu et
117 al., 2025).  We kept only genes with a score of at least five that were also reciprocally
118 the best score between mouse and human and excluded genes with more than one
119 match. This resulted in 16,981 orthologous genes. Cytosolic L and S ribosomal genes
120 were obtained from Human Genome Organization (groups 728 and 729;
121 https://www.genenames.org/data/genegroup/#!/group/). This encompassed 89 human
122 genes, which we subset to the 82 genes with a one-to-one mouse ortholog.
123 Transcription regulator identities were acquired from Animal TFDB (V4; Shen et al.,
124 2023).

### scRNA-seq data acquisition and preprocessing

126 We focused on datasets with count matrices that had cell identifiers readily matched to
127 author-annotated cell types. This was primarily sourced through two means: 1) From the
128 "Cell x Gene" database (https://cellxgene.cziscience.com/), which has pre-processed
129 and annotated data. When a single submission ("collection") contained multiple
130 downloads (for example, different tissue lineages), we downloaded all and combined
131 them into a single dataset keeping only unique cells. 2) Automated screening followed
132 by human curation of the Gene Expression Omnibus (GEO) database (Barrett et al.,
133 2013). Here, we preserved the author-annotated cell types, save for when a biologically
134 uninformative delimiter was used (e.g., "Neuron-1" and "Neuron-2"), in which case we
135 collapsed these cell types into one to prevent overly sparse cell-type populations. We
136 further acquired two tissue-panel datasets. The first was downloaded from the Human
137 Protein Atlas (Uhlén et al., 2015;
138 https://www.proteinatlas.org/download/rna_single_cell_read_count.zip, June 2023),
139 covering 31 tissue-specific datasets which we collapsed into a single dataset and thus
140 treated as a single network. Similarly, we downloaded each of 20 tissue datasets from

141  the Tabula Muris Consortium (2018;
142  https://figshare.com/articles/dataset/Robject_files_for_tissues_processed_by_Seurat/58
143  21263; July 2023), which were also combined as one dataset.

144  Following the advice of the Harvard Chan Bioinformatics Core
145  (https://hbctraining.github.io/scRNA-seq_online/lessons/04_SC_quality_control.html),
146  we uniformly applied relatively lenient filtering rules for all datasets. We required a
147  minimum cell count of 500 UMI (or equivalent) and 250 expressed genes, and a ratio of
148  the $\log_{10}$ count of genes over $\log_{10}$ UMI counts greater than 0.8 for all experiments, save
149  for SMART-seq assays, where the cutoff was relaxed to 0.6 as this technology can
150  result in greater read depth for select genes (Wang et al., 2021). We applied standard
151  CPM library normalization on the raw counts of all datasets (Seurat V4.1.1
152  NormalizeData "RC"), having observed that the log transformation in other normalization
153  schemes resulted in elevated correlation reproducibility in our null comparisons.

154  ***scRNA-seq network construction***

155  Aggregate single cell coexpression networks were constructed as described by Crow et
156  al. (2016). Every dataset consists of a gene by cell normalized counts matrix, where
157  each cell is associated to an annotated cell type. We fix genes to the RefSeq Select
158  protein coding genes, setting unreported genes to counts of 0. This was done so that
159  every resulting network had equal dimensionality.

160  For a given dataset, we performed the following steps for each cell type:

> 161  1. Subset the count matrix to only cells of the current cell type.
> 162  2. Set genes with non-zero counts in fewer than 20 cells to NA.
> 163  3. Calculate the gene-gene Pearson's correlation matrix.
> 164  4. Set NA correlations resulting from NA counts to 0.
> 165  5. Make the correlation matrix triangular to prevent double ranking symmetric
> 166     elements.
> 167  6. Rank the entire correlation matrix jointly, using the minimum ties method.

168  The resulting rank matrices across cell types were then summed into one matrix that
169  was re-ranked and standardized into the range [0, 1] by dividing each element by the
170  maximum rank. Higher values correspond to consistently positive coexpressed gene
171  pairs, and values closer to 0 represent more consistently negative pairs. Step 2 is
172  applied to ensure noisy coexpression values are not calculated from overly sparse
173  populations, as recommended by Ballouz et al. (2015). The zero imputation in Step 4 is
174  to ensure the ranking procedure includes non-measured genes, placing them in
175  between positive and negative correlations. Thus, each dataset is represented as a
176  single gene by gene matrix of coexpression scores aggregated across all labeled cell
177  types. A gene profile refers to a single gene vector (such as a TR gene) from a single
178  matrix; a set of profiles is the collection of profiles extracted from the experiments that
179  measured the given gene.

180  **Gene profile similarity**

181 Coexpression profiles from any one dataset may not have a full complement of
182 measured genes and thus contain tied ranks corresponding to missing values in our
183 framework. Consequently, metrics of similarity that compare all of two lists, such as
184 Spearman's correlation, are inappropriate and so we focused on the agreement of the
185 Top and Bottom K genes between profiles. We calculated various set overlap metrics
186 between lists and, finding our conclusions to be consistent, opted for the interpretability
187 of reporting the size of the $Top_K$ and $Bottom_K$ overlaps. We restrict our reporting to TRs
188 that were measured in at least five datasets.

189 For each TR, we calculated the pairwise similarities among its set of profiles. Averaging
190 these similarity metrics was used to infer the consistency of a TR's coexpression profile
191 across datasets. This process was also applied to each of the 82 ribosomal genes to
192 provide a comparison with a set of genes known to be coexpressed. To generate a null
193 comparison, a random TR was selected from each network to create a set of shuffled
194 profiles, and pairwise similarities were calculated and averaged as above. This process
195 was repeated 1000 times, generating a null distribution of average pairwise similarities.
196 A TR with an average similarity greater than any of the 1000 nulls has an empirical *p-*
197 *value* < 0.001.

**Aggregating TR profiles and the effect of gene measurement sparsity**

199 To prioritize the gene partners most commonly coexpressed with each gene, we
200 averaged the set of rank-standardized profiles for the given gene into one aggregate
201 profile. As each dataset-level profile had variable gene measurement, there was
202 variable delineation between the positive coexpression values, the non-measured gene
203 pair ties, and negative coexpression values. Therefore, for a given gene's set of profiles,
204 we imputed all tied values to the median value before averaging, standardizing the
205 values of non-measured gene pairs. A schematic is shown in Supplemental Fig. 1C.

**Gene set enrichment**

207 For each aggregate profile, we performed GO enrichment analysis of "biological
208 process" terms with the "precRecall" R implementation of ermineJ
209 (https://github.com/PavlidisLab/ermineR; Ballouz et al., 2016), using the aggregate
210 values as scores. This approach considers the full scored list to find enriched terms but
211 places greater emphasis on the top of the gene list. We analyzed 3,284 terms that had
212 20-200 genes and set the false discovery rate at 0.05 for considering terms significant.
213 For the orthologous coexpression rankings we used human genes to map GO
214 annotations.

**ChIP-seq data acquisition and summarization**

216 All ChIP-seq data was downloaded from the Unibind database (Puig et al., 2021;
217 https://unibind.uio.no/downloads/; September 2022). For every TR experiment, we
218 scored gene binding intensity using the same approach as in Morin et al., 2023, using a
219 continuous scoring function (Ouyang et al., 2009; detailed in the Supplement). To
220 generate an aggregate binding profile, we averaged the gene binding vectors specific to

6

221 each TR. A "consensus" list of ASCL1 bound regions consisted of the union of all its
222 peaks across ASCL1 Chip-Seq datasets (detailed in Supplement).

**Literature curation evaluation**

224 TR-target interactions supported by low-throughput experimental evidence were
225 collected from our prior study (Chu et al., 2021), which compiled information from other
226 resources (see Supplement for details) and then significantly expanded upon
227 neurologically-relevant TRs. Since Chu et al. (2021) was published, we have further
228 expanded this collection, to a total of 27,629 experiments encompassing 772 TRs and
229 5,899 gene targets. We then used each TR's aggregate profile's ranking as a score and
230 its curated targets as labels, calculating AUC metrics (AUPRC: area under the precision
231 recall curve and AUROC: area under the receiver operator curve) using the ROCR
232 package (Sing et al., 2005; V1.0-11). To generate a null comparison for each TR, we
233 randomly sampled from the entire literature curation corpus a number of targets equal to
234 the count of curated targets for the given TR, and calculated AUCs using the TR's
235 aggregate profile as a score and the shuffled targets as labels. This process was
236 repeated 1000 times to generate a null distribution of AUC values. The observed AUCs
237 (using the TR's true curated targets) were then compared to this distribution of null
238 AUCs. A quantile of 1 means that the observed AUC was better than every single null
239 AUC (empirical *p-value* < 0.001). We restrict our reporting to TRs that had at least five
240 curated targets.

**Cross-species coexpression profile comparison**

242 There were 1,246 TRs with a one-to-one orthologous match between mouse and
243 human that were also measured in at least 5 datasets in each species. For each of
244 these TRs, we subset their aggregate profiles to the 16, 981 orthologous genes. Each
245 orthologous TR thus has a mouse and human aggregate profile, generated separately
246 across the respective species' datasets. To generate a consensus orthologous profile
247 for each TR, we took the rank product between its human and mouse aggregate
248 profiles. To compare ortholog aggregate profiles, we calculated Spearman's correlation
249 and TopK and BottomK overlaps. Null comparisons were generated in a manner
250 consistent with the individual profile comparison: similarities were calculated between
251 randomly shuffled aggregate profiles between species over 1000 iterations.

252 To learn the specificity of a TR's aggregate coexpression profile with its matched
253 ortholog in the reciprocal species, we combined the framework applied in this study with
254 prior studies examining the conservation of coexpression (Patel et al., 2012; Suresh et
255 al., 2023). For each TR in a species, we selected the given TR's top 200 coexpressed
256 partners ($Top_{200}$). We next calculated the overlap of this gene set with the $Top_{200}$ gene
257 sets of each of the 1,246 TRs in the other species. We then treated the mismatched
258 (non-orthologous) overlaps as a distribution and represented the matched (ortholog)
259 TR's $Top_{200}$ as a quantile with respect to this distribution. We refer to this quantile as the
260 *Ortholog retrieval score*. A score of 1 means that the given TR's ortholog shared more
261 top coexpressed partners than any other TR in the other species. This procedure was
262 then repeated for the reciprocal species. The result is a pair of *Ortholog retrieval scores*
263 for each TR: how well a human TR's aggregate profile recovered its mouse ortholog

7

264  relative to all other mouse TRs (human in mouse), and the recovery of the mouse
265  ranking across human TRs (mouse in human).

266  **Integrating coexpression and ChIP-seq profiles**

267  For TRs with ChIP-seq data, we took the rank product of the TR's aggregate
268  coexpression profile and its aggregate binding profile, re-ranking the result (Breitling et
269  al., 2004; Wang et al., 2013; Morin et al., 2023). This orders genes by placing equal
270  weight on their (positive) coexpression evidence and their binding evidence. We further
271  report a second prioritization scheme for each TR, categorizing genes based on a cut-
272  off of the rankings for both data types and species:

273        1. Stringent: Required a gene's presence in the Top 500 of both coexpression
274        and binding in both species (orthologous genes only).

275        2. Elevated: Genes needed to make the Top 500 cut-off for both data types in
276        one species and in one data type for the other species (orthologous genes only).

277        3. Species-specific: Top 500 cut-off for both data types in one species only.
278        Notably, this may include genes absent from the one-to-one orthologous set, or
279        TRs that had ChIP-seq data in one species only. Consequently, this tier had
280        greater coverage than the others.

281        4. Mixed-species: Allowed genes ranked in the Top 500 in both data types, but
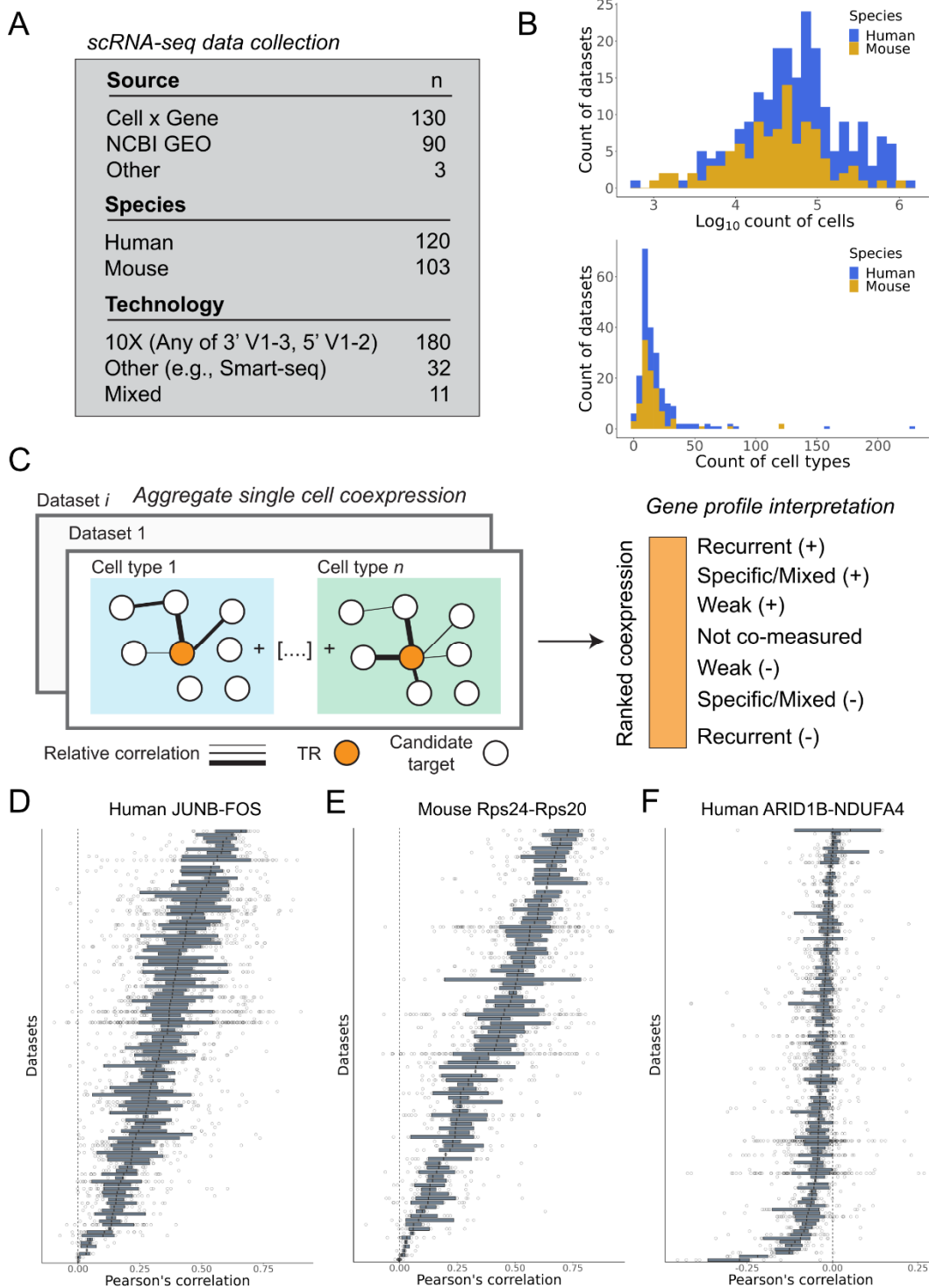282        each in only one species (orthologous genes only).

# Results

### Assembling a broad corpus of single cell RNA-seq data

285  To establish a diverse range of biological contexts for constructing single cell
286  coexpression networks, we acquired scRNA-seq data from public resources (Methods).
287  Our focus was strictly on datasets that included author-annotated cell type labels in the
288  metadata, and all identified datasets underwent consistent preprocessing. In total, we
289  analyzed 120 human datasets and 103 mouse datasets (Fig. 1A; Metadata in
290  Supplemental Table 1). This corpus spans a wide range of biological contexts, scRNA-
291  seq technologies, and counts of assayed cells. After filtering, the median human dataset
292  measured 15,341 protein coding genes across 74,148 cells and 14 cell types; in mouse
293  13,996 genes across 36,755 cells and 12 cell types (Fig. 1B; Supplemental Figs. 1A,B).
294  There was appreciable spread in these counts, with tissue atlas studies typically
295  exhibiting the broadest coverage. The complete dataset is over 2.8 x$10^7$ cells.

### Constructing single cell coexpression networks

297  We constructed aggregated single cell coexpression networks for each dataset using
298  the approach outlined by Crow et al., 2016 (Methods). In brief, this entails generating a
299  gene-by-gene correlation matrix for each cell type within a dataset, ranking each cell
300  type correlation matrix, and consolidating them into a single network per dataset (Fig.
301  1C). Notably, unlike in Harris et al. (2021), where information was consolidated across

**Figure 1.** Overview of study design. (A) Counts of datasets by source, technology, and species. (B) Top panel: Counts of cells across the dataset corpus. Bottom panel: Counts of cell types. (C) Schematic of the single cell coexpression aggregation framework and the interpretation of an individual gene coexpression profile. (D, E) Examples of the most reproducible positively coexpressed gene pairs. Each bar represents a dataset/network, and each point represents the gene pair's correlation in a cell type within the dataset. (F) Example of one of the most reproducibly negative coexpression gene pairs.

310  datasets for a single cell type, we first aggregate across cell types within a dataset
311  before aggregating across datasets. In doing so, we explicitly prioritize signals shared
312  across cell types. This strategy also minimizes effects due to expression differences
313  between cell types, which we consider a separate question from "within cell" regulatory
314  interactions (Farahbod and Pavlidis, 2020).

315  This procedure aims to rank coexpression partners, as illustrated in Fig. 1C, by ordering
316  from "top" to "bottom": consistently high positive interactions across cell types;
317  mixed/specific positive interactions; weak-to-no coexpression; non-measured gene
318  pairs; and then the increasingly most reproducibly negative coexpressed pairs. From
319  this network, it is possible to extract a single gene column (herein, gene profile), such as
320  for a TR, with the relative ordering reflecting the strength of its aggregate transcript
321  covariation with all other genes.

322  While the focus of this study is on TRs, we first examined the globally most consistent
323  coexpressed gene pairs (Figs. 1D-F). Top examples include TRs that dimerize to form
324  the pleiotropic AP-1 complex, such as JUNB and FOS, as well as members of the
325  ribosomal complex. Given the known biological coexpression of ribosomal genes (Li et
326  al., 2016), we use a set of 82 large (L) and small (S) ribosomal genes that are highly
327  conserved between mouse and human as a positive control when examining TR-gene
328  coexpression in the following analyses (Methods). We also show one of the most
329  consistently negative coexpressed TR-gene pairs in human. Aligning with our prior
330  observations (Lee et al., 2004), the magnitudes of these values are smaller and less
331  consistent than the positive coexpression profiles, contributing to the complexity in
332  identifying repressive interactions (discussed in the Supplemental Material).

### *Similarity of TR-target profiles*

334  Before prioritizing reproducible TR-gene interactions, we examined the concordance of
335  the TR coexpression profiles between datasets. We expected that distinct profiles
336  generated for the same TR and similar contexts would have elevated similarity relative
337  to mismatched contexts or gene profiles. At the same time, the underlying data we used
338  was from differing cell types, as datasets could be from different tissues. While we
339  expected this would affect the degree of similarity, a total absence of overlap between
340  profiles would raise questions about the efficacy of our framework in finding
341  reproducible interactions.

342  We report here on the size of the overlap ($K$) of the top positively coexpressed ($Top_K$)
343  genes between each pair of gene profiles (negative coexpression is discussed in the
344  Supplemental Material). We examined a range of $K$, from 200 — approximately the top
345  1% of protein coding genes — to 1000, finding that our main conclusions were robust to
346  this cut-off. To contextualize the similarity between TR profiles, we generated null
347  similarities, iteratively sampling TRs across datasets and calculating the overlap of the
348  shuffled TR profiles. We also report the similarity of the set of 82 L/S ribosomal genes.

349  First, for each TR we pairwise compared its profiles across studies. As expected, the
350  most similar pairs were supported by datasets investigating similar biological contexts.
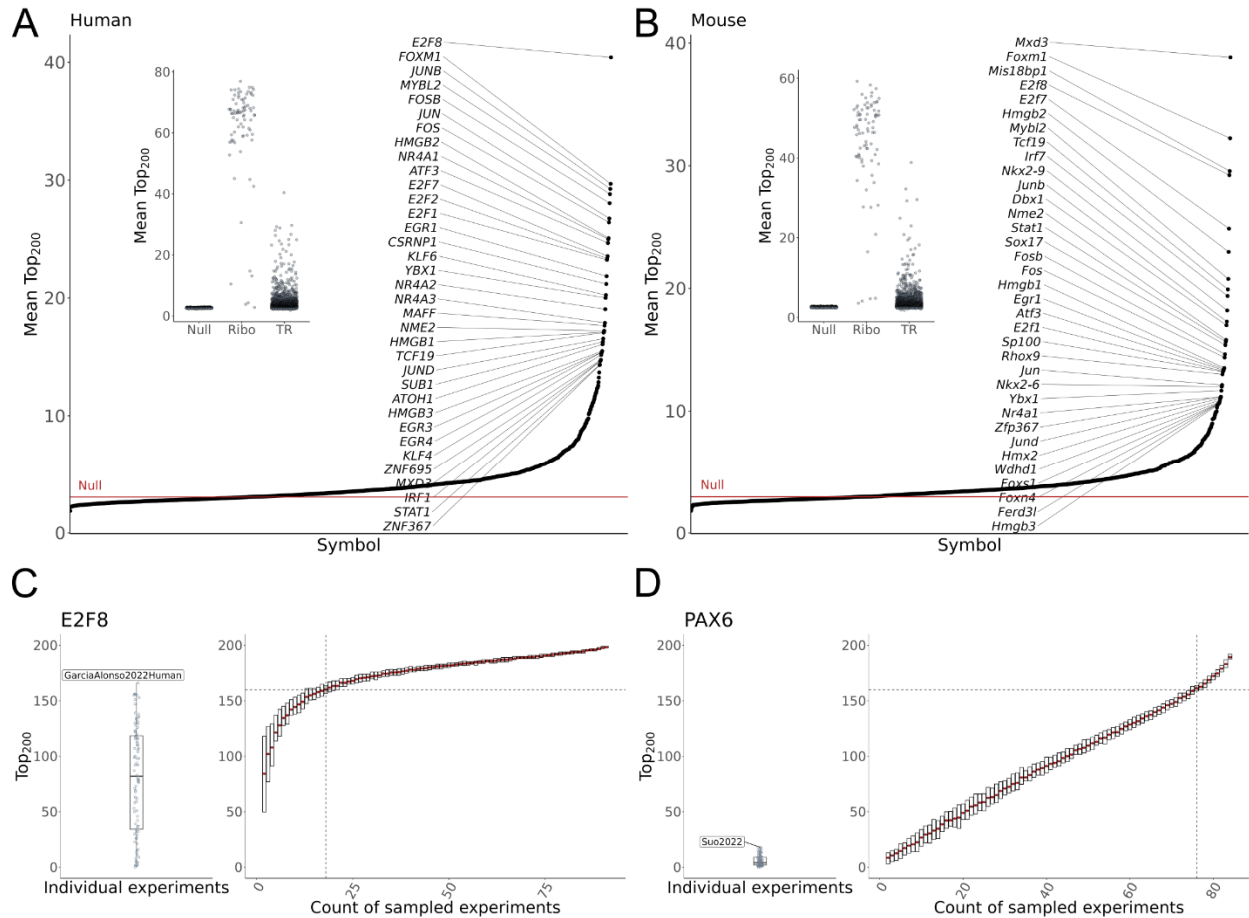351  For example, the best pairing in human ($Top_{200}$ = 163/200) was between *FOXM1*

10

profiles from two studies that both assayed the developing human intestine (Fawkner-Corbett et al., 2021; Elmentaite et al., 2021). The highest Mouse $Top_{200}$ (150/200) was associated with E2f8, derived from two studies of the blood-brain barrier (Posner et al., 2022; van Lengerich et al., 2023). The magnitude of the best ribosomal gene pairs was comparable: the best global human ribosomal pairing ($Top_{200}$ = 161/200) belonged to *RPS13*, originating from two immune cell studies (Liu et al., 2021; Domínguez Conde et al., 2022).

While these observations support the ability to find consistent coexpression patterns within pairs of similar contexts, our ultimate aim was to combine information across contexts. Seeking a more global summary of TR profile overlap, we calculated the mean $Top_{200}$ overlap for each TR profile across all unique pairs of networks measuring the TR. We again use the similarities from the pairs of randomly sampled TRs and the 82 ribosomal genes as reference.

In Figs. 2A,B, we show the average $Top_{200}$ of shuffled TR pairs across 1000 iterations. The typical null sample had an average $Top_{200}$ value of 2.7/200 in human and 2.6/200 in mouse. The ribosomal genes, approximating an empirical "upper bound," averaged 61/200 in human and 44/200 in mouse. The distribution of average $Top_{200}$ values was highly skewed for TRs, with 67% of human TRs and 68% of mouse TRs having an average $Top_{200}$ value greater than the maximum value achieved across all of the null samples (empirical *p-value* < 0.001; represented as red lines in Figs 2A, B). And while the best individual ribosomal data pairs were equivalent in overlap size compared to the best individual TR pairs, ribosomal genes typically had a much greater average $Top_{200}$ than even the best TR. This underscores the unusual uniformity of ribosomal protein gene coexpression across distinct cellular contexts — it is an outlier. A similar comparison for the $Bottom_{200}$ is provided in Supplemental Figs. 2A-D.

TRs with the highest mean $Top_{200}$ values, indicative of the most consistent positive coexpression profiles across studies, were often associated with fundamental cellular housekeeping processes. For example, *E2F8* led in human (mean $Top_{200}$ 40.4/200), with mouse *E2f8* similarly having one of the most consistent profiles (Figs. 2A,B). The E2F family are well characterized regulators of the cell cycle (Emanuele et al., 2020), and other E2F members also ranked high in both species, as did regulators involved in early transcriptional response to environmental signals, such as AP-1 complex members *FOS* and *JUN*. In mouse, the highest mean $Top_{200}$ belonged to *Mxd3*, a MYC-antagonist whose human ortholog also had elevated similarity. More broadly, there was appreciable correlation between human and mouse orthologous TRs (Methods) in the consistency of their positive and negative coexpression profiles (Supplemental Figs. 3A,B).

TRs with context-restricted activity might be expected to exhibit relatively low cross-dataset similarity in our broad corpus. However, this is not necessarily the case. For example, the neural regulator NEUROD6 (Tutukova et al., 2021) had one of the most consistent TR profiles in human (mean $Top_{200}$ rank 44th out of 1,605 TRs), despite being only measured in 22 of 120 datasets. This shows that restricted expression does not preclude the identification of reproducible patterns. In contrast, human PAX6 —

11

**Figure 2.** Similarity of TR profiles. (A) Inset: distribution of the mean $Top_{200}$ overlaps for the null background, 82 ribosomal genes, and 1,605 human TRs. The null was generated through 1000 iterations of sampling one TR profile from each of 120 human datasets and calculating the average size of the $Top_{200}$ overlap between every pair of sampled profiles. The ribosomal genes represent a "base case" scenario. Main: The average $Top_{200}$ overlap of all human TRs, with the red line indicating the best null overlap. (B) Same as in A, save for 103 mouse experiments and 1,484 TRs. (C,D) Saturation analysis of global TR profiles for human (C) E2F8 and (D) PAX6. Left panels show the spread of $Top_{200}$ overlaps between individual dataset profiles and the global E2F8 and PAX6 profiles. Right panels show the spread of overlaps when iteratively subsampling and aggregating datasets at increasing steps. Dotted lines indicate the average number of sampled datasets required to reach 80% of the global profile. E2F8 recovers its global profile with relatively fewer datasets than does PAX6.

12

409  necessary for the development and function of several nervous and pancreatic tissues
410  (Wen et al., 2009; Yeung et al., 2016) — had a mean $Top_{200}$ value marginally above the
411  null, improving slightly at K=1000 (Supplemental Figs. 2E,F). Although PAX6 can also
412  be described as a context-restricted, it was detected in 85 of 120 datasets, suggesting
413  greater heterogeneity in its coexpression profiles compared to NEUROD6.

### Ranking aggregated coexpression to prioritize TR-target candidates

415  The preceding section demonstrated that similar TR profiles could be identified across
416  this biologically heterogeneous corpus, supporting the potential to find reproducibly
417  coexpressed gene pairs. We thus turned to our primary aim of prioritizing these
418  consistent interactions, generating a unified gene ranking for each TR using all
419  compiled data. This process involves aggregating information at two levels: first, across
420  cell types *within* a dataset (as in the previous section; Fig. 1C), and then, for each TR,
421  aggregating their profiles *across* datasets (Methods, Supplemental Fig. 1C). This
422  approach aims to maintain the interpretability of an aggregate profile relative to a profile
423  from an individual network (Fig. 1C): the extremes represent the most consistent
424  positive and negative correlations, while the middle of the list encompasses weak and
425  non-measured coexpression gene pairs.

426  As before, we used the set of ribosomal genes to validate that our aggregation workflow
427  prioritized known biological coexpression (Supplemental Material; Supplemental Figs.
428  4A,B). We next performed GO biological process enrichment on all aggregate profiles
429  (Supplemental Fig. 4C), finding that most TRs (91% human, 86% mouse) were
430  associated with at least one term (FDR 0.05). E2F8 coexpression partners were
431  enriched for multiple terms relating to cytokinesis and chromosomal organization, as
432  expected for its known role in these processes (Emanuele et al., 2020). We also
433  frequently observed that terms affiliated with tissue-specific processes were enriched for
434  TRs implicated in those tissues. Examples include glial development and myelination
435  terms for the oligodendrocyte TRs OLIG1/2 (Szu et al., 2021), neuronal synaptic
436  functionality for the aforementioned NEUROD6 (Tutukova et al., 2021), leukocyte and
437  cytokine processes for IRF8 (Salem et al., 2020), and hematopoietic terms for the
438  erythroid GATA1 (Ferreira et al., 2005). Some tissue-selective TRs were enriched for
439  more general regulator terms (e.g., "cell fate commitment" for mouse Pax6) or had
440  disparate tissue-specific terms (e.g., "regulation of osteoblast differentiation" and
441  "regulation of neuron differentiation" for SOX4), potentially reflecting data heterogeneity.
442  While GO is an imperfect resource, these results agree with our other observations that
443  our analysis yields biologically-relevant signals.

444  We examined the relationship between the aggregated global TR profiles and the
445  constituent datasets through two analyses. First, we assessed how well individual
446  experiments aligned with the global profiles to identify potential biases (Supplemental
447  Material). As shown in Supplemental Figs. 2H-L, datasets with the highest agreement
448  were large studies of broad tissues using the 10X Chromium platform, though
449  consistencies between platforms were still observed (Supplemental Fig. 2G).

450  Second, we performed a saturation analysis to determine how many datasets are
451  needed, on average, to recover each TR's global profile (≥80% overlap in $Top_{200}$

452  genes). By iteratively subsampling and aggregating datasets, we evaluated the
453  convergence of sampled TR profiles to the global set. For example, E2F8 (Fig. 2C)
454  required an average of 18 of 92 datasets to reach saturation, while PAX6 (Fig. 2D)
455  showed a linear trend, indicating saturation has not yet been achieved. These results
456  suggest future work is needed to explore not only replicable context-specific patterns for
457  TRs such as PAX6, but also the extent to which globally consistent partners can be
458  found when using more data.

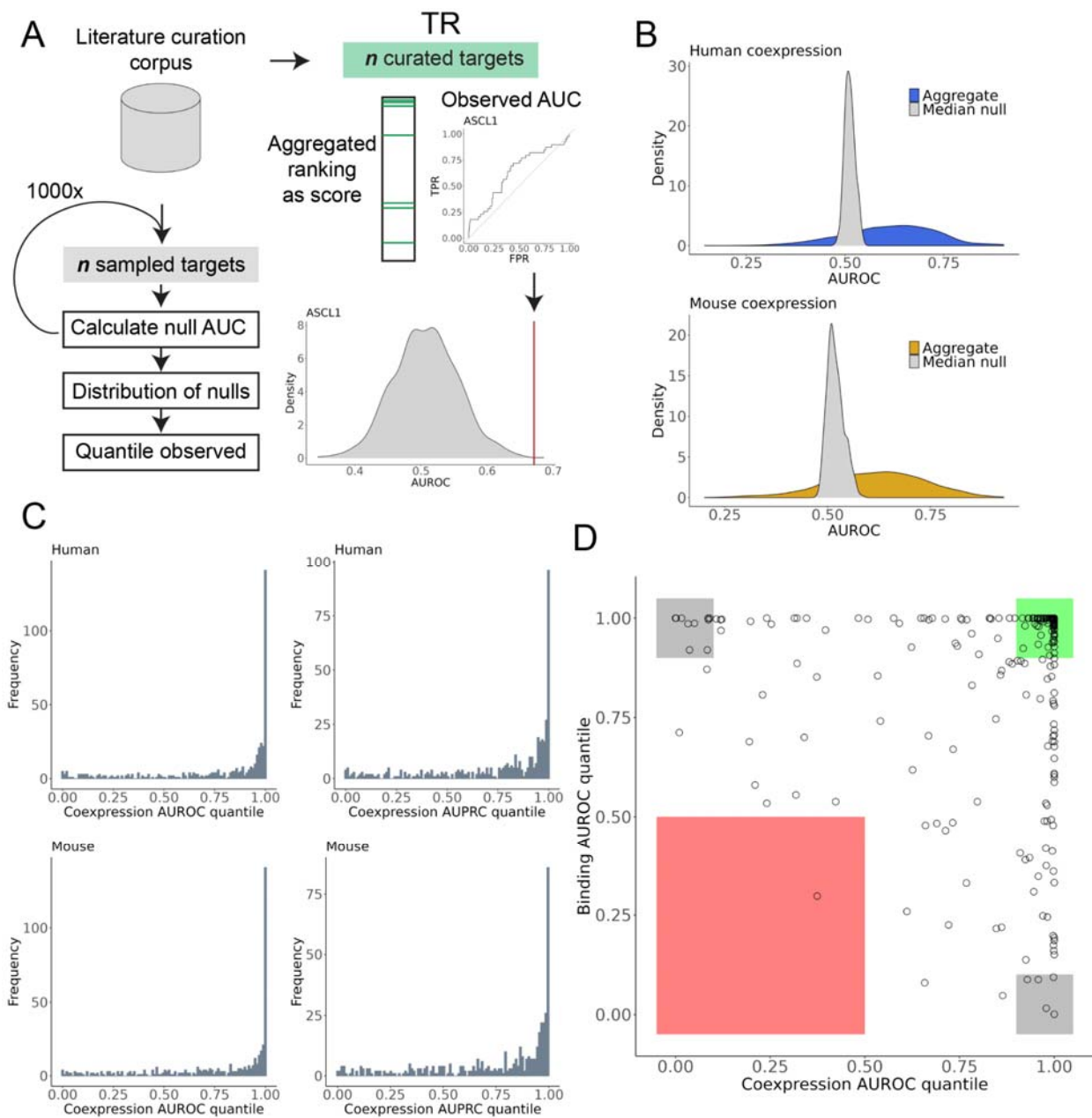### Recovery of literature-curated TR-target interactions

460  Equipped with a unified single cell coexpression profile for each human and mouse TR,
461  we aimed to assess their concordance with an orthogonal line of regulation evidence.
462  While coexpression is expected to prioritize both direct and indirect regulatory
463  interactions (the latter we would consider false positives), the rankings should still
464  demonstrate a greater ability to recover true direct interactions relative to a null
465  expectation.

466  In a previous study (Morin et al., 2023), we evaluated the utility of aggregating TR
467  perturbation and ChIP-seq experiments, using literature-curated low-throughput
468  interactions as positive labels and calculating area under the curve (AUC) metrics
469  (Marbach et al., 2012; Garcia-Alonso et al., 2019). We applied the same framework
470  here, using curated TR-target interactions we have collected (Chu et al., 2021, since
471  expanded) and assembled from other resources (see Supplement for further
472  discussion). We considered TRs that had a minimum of five curated targets, resulting in
473  451 TRs analyzed in human (median count of curated targets = 18) and 434 in mouse
474  (median count = 17).

475  We first examined the effectiveness of the aggregate profiles in recovering curated
476  targets relative to the individual TR profiles that compose each aggregate. On average,
477  the aggregate profiles outperformed (better prioritized curated targets) the expected
478  AUC value from an individual profile (Supplemental Fig. S5A). Therefore, aggregating
479  the coexpression networks typically maintains or improves performance on this
480  benchmark.

481  Next, we evaluated the efficacy of the coexpression rankings in recovering curated
482  targets relative to a null distribution of AUCs (*Quant_coexpression)*. While the raw AUC
483  values were typically better than random (Fig. 3B, Supplemental Fig. 5B), we report the
484  quantile of the observed value relative to a null to standardize the comparison across
485  TRs (discussed in Supplemental Material). This null was created by size-matching and
486  randomly sampling from the pool of curated targets from the entire literature-curation
487  corpus. The latter helps account for biases in the coverage of targets in the low-
488  throughput literature. A *Quant_coexpression* value of 1 indicates that an aggregate
489  profile outperformed every null sample.

490  ASCL1 is provided as an example of this procedure for one TR in Fig. 3A. As illustrated
491  in Fig. 3C, the coexpression aggregates consistently exceeded the null AUCs, reflected
492  by a median AUROC *Quant_coexpression* of 0.95 in human and 0.93 in mouse. The
493  pile-up of quantiles near or equal to 1 indicates that, while not universal, a majority of

**Figure 3.** Recovery of literature curated targets by aggregate rankings. (A) Schematic of literature curation evaluation. (B) Distributions of the observed AUROCs for 451 human and 434 mouse aggregate TR coexpression profiles, along with the distribution of the median null AUROCs generated for each profile. (C) Histograms of the AUROC and AUPRC coexpression quantiles for human and mouse. (D) Scatter plot of the AUROC quantiles for the coexpression and binding profiles of 253 human TRs that had binding data and at least five curated targets. Green box indicates TRs for which both genomic methods were effective in the benchmark, grey box for only one method, and red box for neither method being effective.

503    TR single cell coexpression rankings excelled in prioritizing matched curated targets
504    over randomly sampled targets. These observations strongly suggest that these
505    aggregate rankings are capable of prioritizing regulatory interactions that were identified
506    through targeted biochemical assays.

507    To further contextualize these performances, we conducted a similar null AUC analysis,
508    this time using aggregate ChIP-seq signals. In brief, we applied the same approach as
509    in Morin et al., 2023, scoring gene-level binding intensity for each ChIP-seq experiment,
510    then averaging these signals within each TR's set of experiments to create a single
511    unified ranking of gene binding for each TR. In total, we considered 4,115 human
512    experiments for 253 TRs and 3,564 mouse experiments for 241 TRs from the Unibind
513    database (Puig et al., 2021, Methods) that had at least five curated targets. As with the
514    aggregate coexpression signal, we compared the unified binding ranking's ability to
515    recover TR-specific curated targets relative to a null of sampled targets
516    (*Quant_binding*).

517    We anticipated that TR ChIP-seq, as a more direct form of regulatory inference, might
518    outperform coexpression (Garcia-Alonso et al., 2019). However, in our hands the
519    aggregate binding evidence was on par with single cell coexpression in its ability to
520    predict known targets (Supplemental Fig. 5C), further motivating integration of both data
521    types. Supporting this, integrating the coexpression and binding rankings for available
522    TRs typically led to elevated performance in the benchmark (Supplemental Fig. 5D).

523    Among TRs with both binding and coexpression data, many performed well in the
524    benchmark for both data types separately, as demonstrated for human TRs in Fig. 3D.
525    In human, 134 of 253 (53%) TRs had AUCs (AUPRC or AUROC) *Quant_binding* > 0.9
526    and *Quant_coexpression* > 0.9; in mouse 126 of 241 (52%). This signifies that, for these
527    specific regulators, aggregated coexpression and binding profiles both effectively
528    prioritize curated TR targets relative to sampled targets. This alignment highlights TRs
529    whose activity may be more readily identified through distinct data modalities. Further,
530    of the TRs performant in both lines of evidence, more than half did so in both species
531    (human 83 of 134, mouse 83 of 126), suggesting convergence of evidence across not
532    only experiments, but also species.

533    This agreement of evidence encompassed broadly active TRs, such as those involved
534    in the AP-1 complex. However, it also included more specialized factors, such as the
535    neuronal-specifying ASCL1, and the aforementioned PAX6. This suggests that, even
536    though the average overlap of PAX6 profiles was weak (Fig. 2D; Supplemental Figs.
537    2E,F), there was still a consensus of recurrent curated PAX6 targets within these
538    smaller intersects. We also find cases where only one data type was performant. LEF1,
539    for example, had an AUROC *Quant_coexpression* value of 1 in both species but a
540    *Quant_binding* value of 0 and 0.22 in human and mouse, respectively.

541    Finally, because negative expression correlations might be of interest for identifying
542    repressive interactions, we conducted an analysis of the reproducibility and
543    performance of the relations predicted from the bottom of the rankings. We found that
544    for some TRs, negative correlations performed better than positive correlations in the
545    benchmark, though this was the exception (Supplemental Fig. 5B). This suggests that

16

546  for some TRs, repressive activity might be inferable from coexpression (see
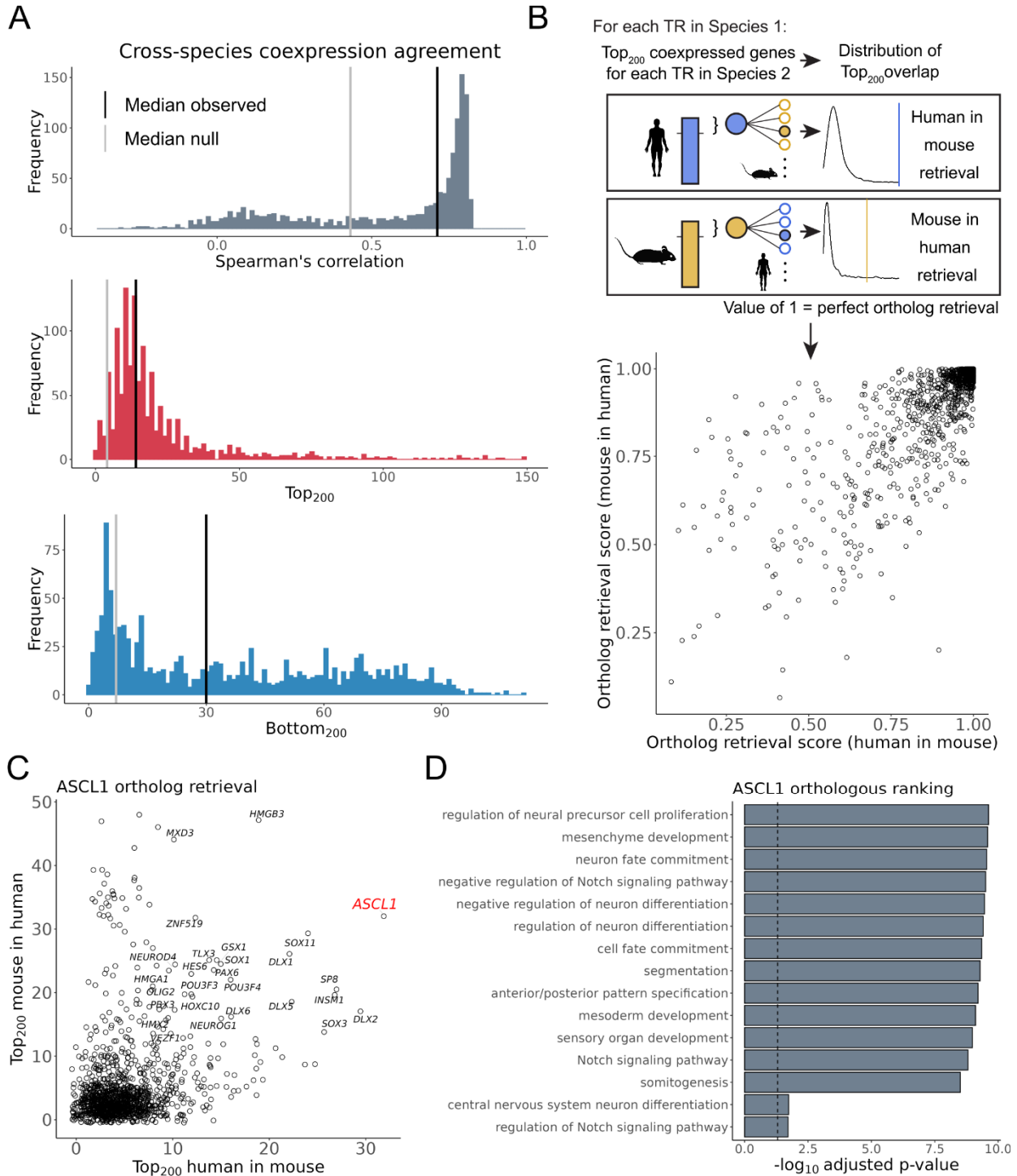547  Supplement for discussion; Supplemental Tables 2-3).

548  ***Identification of conserved interactions***

549  It has been observed that, despite the high evolutionary turnover of regulatory DNA
550  sequences, TR-target relations exhibit relatively high conservation (Yue et al., 2014),
551  with coexpression providing an attractive means to nominate common and divergent
552  interactions (Monaco et al., 2015; Lee et al., 2020; Suresh et al., 2023). Here, our aim
553  was to identify the extent to which individual TR aggregate coexpression profiles were
554  preserved between mouse and human, focusing on orthologous genes (Methods). A
555  meta-analytic comparison of TR single cell coexpression profiles between these two
556  species is lacking, and we reasoned that evidence of conservation using this global data
557  corpus would provide future support for studies that focus on specific TR patterns in a
558  more focused context.

559  Figure 4A demonstrates the similarity distributions between ortholog aggregate
560  coexpression profiles, overlaid with the median observed and shuffled null values.
561  Although there was appreciable spread in these similarity metrics, most orthologs
562  shared more similarity in their profiles than would be expected from shuffled TRs,
563  suggestive of conserved TR coexpression. While there are TRs that agree less well
564  between species, we are cautious in interpreting this as species-specific regulatory
565  rewiring, given the relatively modest effect size and the absence of an exact match in
566  cellular contexts covered across both species.

567  Given our emphasis on reproducible interactions, we focused on the overlap at the
568  extremes of these species rankings (Figs. 4B,C; Supplemental Fig. 6C). To quantify the
569  specificity of this overlap, we applied a slightly modified framework of the $Top_K$ overlap
570  used in this study, consistent with prior studies (Methods; Patel et al., 2012; Suresh et
571  al., 2023) and illustrated in Fig. 4B. The result is a pair of ortholog retrieval scores for
572  each TR: how well a human TR's ranking recovered its mouse ortholog relative to all
573  other mouse TRs (human in mouse), and the recovery of the mouse ranking across
574  human TRs (mouse in human), with a value of 1 indicating perfect retrieval.

575  As demonstrated in Fig. 4C, there was considerable preservation of single cell
576  aggregate TR coexpression profiles between mouse and human. The median ortholog
577  retrieval score for human was 0.969, with 175/1,246 (14%) TRs having a perfect value
578  of 1; in mouse these values were 0.973 and 172/1,246 (14%), respectively. These
579  relative values correspond to a median $Top_{200}$ overlap of 14 genes, with FOXM1 and
580  HMGB2 each having a maximal $Top_{200}$ of 149 genes (Fig. 4A). While the most
581  conserved TRs (by $Top_{200}$ overlap) were led by regulators of housekeeping processes
582  such as cell division, we also observed this preservation among more specific TRs,
583  such as the aforementioned NEUROD6 (human in mouse and mouse in human = 1,
584  $Top_{200}$ = 50). Logically, many of these highly preserved TRs also had similar profiles
585  across datasets within species (as shown in Figs. 2A,B), and those that were weakly
586  preserved generally lacked consistency within species (Supplemental Figs. 6A,B).
587  These findings collectively contribute to characterizing the extent to which each TR can

17

588
589 **Figure 4.** Preservation of mouse and human single cell coexpression profiles. (A) Distribution of
590 coexpression agreement between the aggregate single cell coexpression profiles of 1,246
591 orthologous TRs. Black lines indicate the median value for the TRs, grey lines indicate the
592 median of null values generated by shuffling pairs of orthologous TRs. (B) Top: Schematic of
593 the ortholog retrieval workflow, adapted from Suresh et al., 2023. Bottom: Scatterplot of the
594 resulting ortholog retrieval scores (C) Scatter plot of the ASCL1 $Top_{200}$ overlaps. (D) The top 15
595 GO terms when combining the human and mouse top ASCL1 coexpressed gene partners.
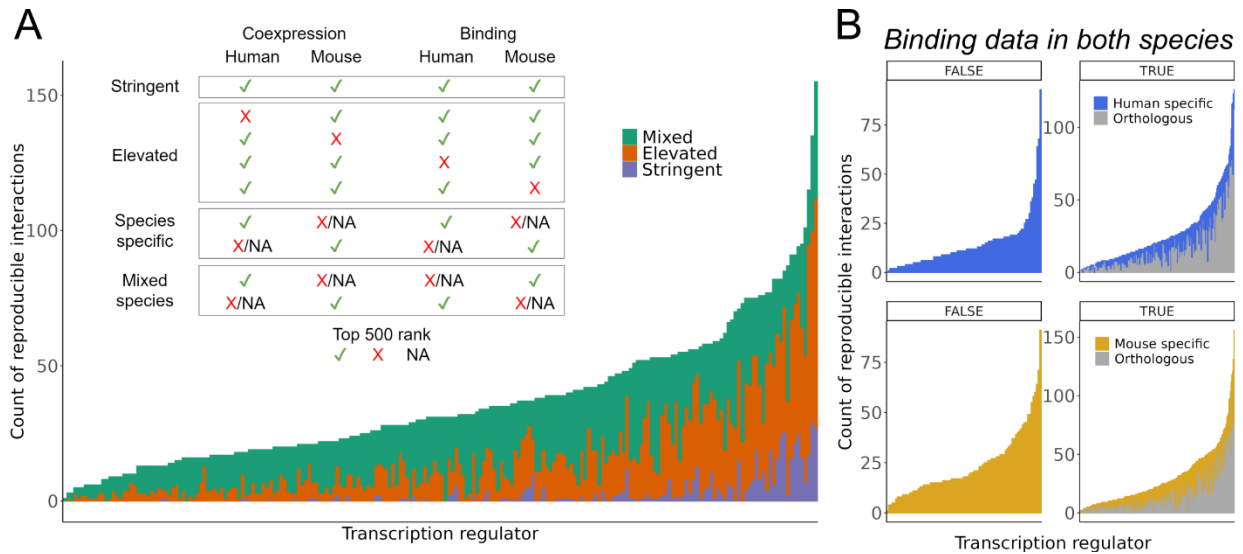
596  be defined by a set of coexpressed gene partners, facilitating inferences into their
597  biological roles.

598  In Fig. 4C we illustrate this overlap procedure for ASCL1, an essential pioneer nervous
599  system regulator that is also relevant to cancer (Castro et al., 2011). Of the 200 genes
600  that were most consistently coexpressed with human ASCL1, 32 of their mouse
601  orthologs were also in the mouse Ascl1 Top$_{200}$ set. This marked the largest overlap
602  human ASCL1 had with any mouse TR (human in mouse = 1). In the reciprocal
603  comparison, where mouse Ascl1 was queried against all human TRs, human ASCL1
604  ranked 30th (mouse in human = 0.98). The 29 human TRs with a greater overlap with
605  mouse Ascl1 did not have a sizable overlap in the reciprocal comparison, save for
606  HMGB3. Conversely, TRs other than ASCL1 with elevated overlap across species
607  included the ASCL1 curated targets INSM1, HES6, and DLX5 (Castro et al., 2006;
608  Nelson et al., 2009; Kito-Shingaki et al., 2014). Other TRs are well-characterized for
609  operating in a regulatory network with ASCL1 — though not necessarily as direct
610  downstream targets — such as DLX1/2/6, GSX1/2, SP8, and OLIG2 (Wang et al., 2013;
611  Al-Jaberi et al., 2015; Liu et al., 2017; Aslanpour et al., Lunden et al., 2019; 2020). GO
612  enrichment of the top ASCL1 coexpressed gene partners using information from both
613  species returned numerous terms that are consistent with ASCL1's role in brain
614  development (Fig. 4D).

615  ***Combining single cell coexpression and aggregated binding reveals numerous***
616  ***reproducible interactions***

617  Up to this point, we have presented evidence supporting the existence of recurrent
618  single cell TR-gene coexpression patterns within (Fig. 2) and across species (Fig. 4),
619  demonstrating that this information can prioritize curated experimental interactions (Fig.
620  3). One of our primary motivations is to prioritize the direct gene targets of TRs (Morin et
621  al., 2023). However, the correlation of TR-gene transcripts serves as an indirect form of
622  gene regulation evidence — it does not confer information about the causative
623  directionality of this covariation. We thus now turn to identifying interactions
624  corroborated by TR binding evidence, using the same aggregated Unibind ChIP-seq
625  data examined in the literature curation evaluation. We reasoned that, as in our earlier
626  work, knowledge of binding can help focus attention on expression patterns more likely
627  to reflect direct regulatory relations.

628  We present two straightforward strategies for prioritizing reproducible interactions,
629  acknowledging the use of relatively arbitrary cut-offs for the sake of reporting. All
630  summarized rankings are made available for researchers interested in conducting their
631  own exploration. We first combined the single cell coexpression and binding profiles into
632  a final ordered ranking for TRs with ChIP-seq data, using the common rank product
633  summary (Breitling et al., 2004; Wang et al., 2013; Morin et al., 2023). This was done
634  separately for each species (317 TRs in human, 305 in mouse), as well as across
635  species for orthologous TRs with available data (216 TRs). This establishes convenient
636  lists that order the protein coding genes most associated with each TR based on their
637  aggregated single cell coexpression and binding profiles.

19

**Figure 5.** Count of interactions supported across methods and species. (A) Inset: criteria used to group interactions into tiers. Bar chart: Count of unique interactions gained in each orthologous tier (Stringent, Elevated, and Mixed-Species) for the 216 TRs with binding data in both species. (B) Count of Species-Specific interactions for 317 TRs in human (top) and 305 TRs in mouse (bottom). TRs are split by those with ChIP-seq data in one species only (left) and thus are ineligible for consideration in the orthologous interactions, and those with ChIP-seq data in both species (right). Grey bars indicate the count of interactions already found in the Stringent and Elevated sets, coloured bars indicate the count of Species-Specific interactions that were gained due to lacking orthologs or because they had elevated ChIP-seq signal in one species and not the other.

649    Recognizing that a gene may be prioritized (have a better rank product) if ranked
650    exceptionally well in one data type or species only, we introduce a second scheme for
651    more balanced consideration across lines of evidence. For each TR, genes are
652    categorized into tiers by their status across the rankings, as illustrated in the inset of
653    Figure 5A. This collection provides examples of regulatory interactions supported by
654    both binding and single cell coexpression evidence.

655    Fig. 5A shows the counts of unique orthologous interactions gained in each tier of
656    evidence for the available TRs. The Stringent level, representing the most reproducible
657    interactions across both species and genomic methods, contains 545 TR-gene pairs
658    corresponding to 101 TRs and 357 unique genes. The TRs with the largest Stringent
659    collection featured multiple AP-1 members, led by FOSL1 with 29 genes, along with
660    immunity TRs such as STAT1, STAT2, and IRF1. More specialized TRs also had
661    among the largest Stringent sets, such as the hematopoietic factors SPI1 (n = 27),
662    GATA1 (n = 16) and GATA2 (n = 11), and the hepatic HNF4A (n = 8). This once again
663    suggests conservation of many regulatory interactions, although it is essential to note
664    that this observation is influenced by the limited coverage of ChIP-seq data across
665    biological contexts.

666    The Elevated collection relaxes the criteria to allow orthologous genes reaching the cut-
667    off in three of the four rankings. This resulted in 3,106 Elevated TR-gene pairs, with 211
668    of the 216 available TRs having at least one gene in their set (median = 10). TRs with
669    the largest Elevated collection closely overlapped with those having the largest
670    Stringent sets, reinforcing the notion of preserved target genes among these TRs. The
671    Species-specific level encompasses two groups of TRs: those that have ChIP-seq data
672    in both species and those in only one. This is reflected in Fig. 5B, where we show the
673    count of reproducible interactions for each group. The left panels display TRs with ChIP-
674    seq in only one species and were thus ineligible for consideration in the Stringent or
675    Elevated tiers. In human, this corresponded to 99 TRs with a median of 11 interactions.
676    TFDP1 led with 93 genes supported by both aggregated single cell coexpression and
677    binding evidence. In mouse, all 89 available TRs were associated with at least one gene
678    (median = 18), with the interferon TR Irf8 having a maximum of 91 genes, including
679    numerous immunity-associated genes such as *Mpeg1*, *Ctss*, *Cd180*, *Xcr1*, and
680    *Trim30a*.

### *Highlighting ASCL1*

682    We conclude by focusing on ASCL1, emphasizing that this exploration of ASCL1
683    regulatory targets is just one example made possible by the information we have
684    summarized and made available for community use. In Fig. 6A we present the genes in
685    each tier of evidence for ASCL1, along with their curation status from the 39 available
686    ASCL1 targets in the literature corpus. Human ASCL1 was measured in 61 of 120
687    scRNA-seq datasets, and in mouse 65 of 103. Regarding ASCL1 binding data, there
688    were 10 ChIP-seq datasets in human — largely in cancer cell lines — as well as 10 in
689    mouse, mostly in neuronal and embryonic contexts.
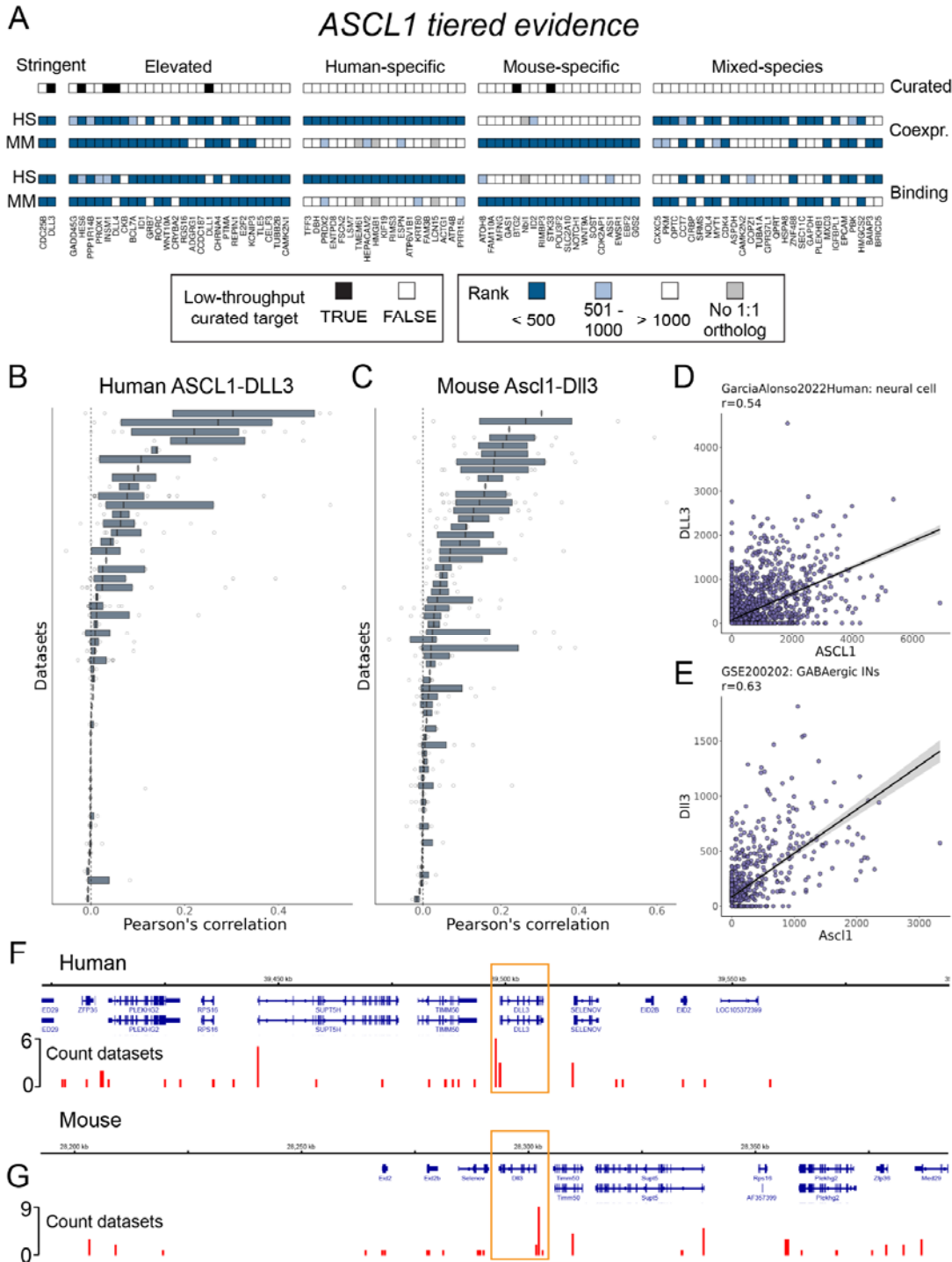
690    Two genes fit the Stringent criteria used for this report: the literature-curated ASCL1
691    target and Notch signalling ligand DLL3 (Henke et al., 2009), and the cell cycling

692    phosphatase CDC25B, which was not in the low-throughput literature collection but is
693    nevertheless discussed elsewhere as a target of ASCL1 (Castro et al., 2006). The
694    Elevated set consisted of 26 genes, with 6 narrowly missing the Stringent criteria
695    (indicated by lighter shading in Figure 6A). Among them are well-described and
696    literature-curated ASCL1 targets, such as the Notch effector HES6 (Nelson et al., 2009)
697    and the neuroendocrine regulator INSM1 (Jacob et al., 2009; Jia et al., 2015). ASCL1
698    and INSM1 serve as markers for neuroendocrine tumours, such as for small cell lung
699    carcinoma (SCLC; Zhong et al., 2022). Another Elevated ASCL1 gene, CKB, has
700    upregulated expression in both SCLC (Borromeo et al., 2016; Qu et al., 2022) and
701    ASCL1-high atypical teratoid/rhabdoid tumours (Tamrazi et al., 2019), suggesting an
702    ASCL1 interaction with oncogenic potential across various contexts. We additionally
703    draw attention to the BAF chromatin remodeler BCL7A, for which we found no ASCL1
704    connection in the literature, and which is also associated with diverse cancers (Baliñas-
705    Gavira et al., 2020; Liu et al., 2021).

706    Other Elevated interactions help characterize ASCL1 as a regulator of both neuronal
707    and oligodendrocyte lineages. This includes the cell cycle regulator GADD45G (Huang
708    et al., 2010), the neuronal tubulin TUBB2B (Mazurier et al., 2014; Lin et al., 2017), and
709    acetylcholine receptor subunit CHRNA4 (Ueno et al., 2012). PPP1R14B and ASCL1
710    expression was used to define a primitive oligodendrocyte progenitor population (Weng
711    et al., 2019). We were unable to find (from a low-throughput study or otherwise) a direct
712    connection between ASCL1 and the neuronal adhesion ADGRG1 (Simão et al., 2018),
713    the cortical-marker and calcium-binding regulator KCNIP3 (Ragazzini et al., 2023), or
714    the neuronal splicing factor CELF3 (Yu et al., 2017), although the latter is used as a
715    neuroendocrine marker to characterize ASCL1-high SCLC subtypes (Zhang et al.,
716    2018). Finally, we highlight REPIN1, an Elevated gene that lacked any ASCL1
717    connection in the literature that is also generally understudied.

718    The next tier, of Species-Specific sets, each comprised 19 genes. PRDX2, for example,
719    is a neuronal-enriched mitochondrial gene that has been shown to enhance ASCL1-
720    induced astrocyte-to-neuron reprogramming (Russo et al., 2021). HEPACAM2 is
721    another gene implicated in cancer (Deprez et al., 2020; Yamada et al., 2022) that we
722    could not find a direct ASCL1 association in the literature. TMEM61, lacking a 1:1
723    mouse ortholog, was only eligible for consideration in the Human-specific set, while the
724    reciprocal applied to the mouse Nbl1. Of the 27 genes in the final tier, the Mixed-
725    Species set, we highlight CXXC5. This zinc finger TR was initially characterized as a
726    bone morphogenic-responsive regulator of Wnt signaling in neural stem cells
727    (Andersson et al., 2009), and has been further described as a signal integrator in
728    development and homeostasis with tumour suppressive qualities (Xiong et al., 2019).
729    These examples collectively illustrate the diverse roles of essential TRs, such as
730    ASCL1, in development and disease.

731    Lastly, we summarize the compiled evidence for the Notch ligand encoding *DLL3*, a
732    well-established and curated ASCL1 target (Henke et al., 2009) that was present in the
733    Stringent collection. DLL3 ranked fourth in the ASCL1 coexpression rankings in both
734    species, making it one of ASCL1's most reproducible coexpression partners. Figs. 6B,C
735    illustrates the distribution of Pearson's correlations for the 238 annotated cell types from

**Figure 6.** Reproducible ASCL1 interactions. (A) Heatmap representing the tiered evidence for ASCL1 candidate targets. (B, C) Distribution of Pearson's correlations for ASCL1-DLL3 in (B) human and (C) mouse, as in Fig. 1E-G. (D, E) Scatterplot of the CPM values for ASCL1 and DLL3 for the cells belonging to the cell type that had the highest correlation in the entire corpus for (D) human and (E) mouse. (F, G) Genome track plots centered on DLL3 (yellow boxes) in (F) human and (G) mouse, where the base of the red bars indicates ASCL1 binding regions, and the height indicates the count of ASCL1 ChIP-seq datasets with a peak in the region.

23

744    54 human datasets in which ASCL1 and DLL3 were co-measured (275 cell types in 61

745    datasets for mouse). Notably, despite being one of the most reproducible ASCL1

746    coexpressions, this association is not universal across all cell types. Figs. 6D,E shows

747    the scatter plots of the individual cell types in which the greatest correlation was found:

748    in human, annotated as "neural cells" ($r = 0.54$; Garcia-Alonso et al., 2022), and in

749    mouse, "GABAergic INs" (interneurons) ($r = 0.63$, Hamed et al., 2022). Given the

750    importance of ASCL1 regulation of Notch signalling in neuronal cells (Castro et al.,

751    2006; Castro et al., 2011; Lampada and Taylor, 2023), these collective observations

752    support that our resource can still prioritize specific interactions.

753    In Figs. 6F,G, we demonstrate the ASCL1-DLL3 binding evidence; DLL3 was ranked

754    493rd in the human aggregate binding profile and 81st in mouse. In human, this

755    corresponded to 83 discrete bound regions (Methods) within 500Kb of either direction of

756    the DLL3 TSS, and 25 within 100Kb; in mouse 73 regions within 500Kb and also 25

757    within 100Kb. We calculated which regions were most frequently bound by ASCL1

758    across datasets, reasoning that this may help prioritize functional ASCL1-DLL3

759    enhancers (while being cognizant of biasing factors like open promoters). Using the

760    500Kb cut-off in human, we found that 20 sites were bound in more than one dataset,

761    and that a region approximately 775 base pairs upstream of the DLL3 TSS had a

762    maximum count of 6. In mouse, 28 regions were bound across multiple datasets, with

763    the most frequently bound region (nine of ten datasets) occurring approximately 400

764    base pairs upstream of the DLL3 TSS.

# Discussion

766    In this study we pursued two main objectives. First, we aimed to understand the

767    behavior of the meta-analytic strategy of aggregating single cell coexpression networks

768    (Crow et al., 2016), applying this methodology across a large and broad corpus of

769    scRNA-seq studies. We believe this technique holds great potential in uncovering

770    robust gene coexpression patterns free from the confounding effect of cellular

771    composition. However, before considering specific cell types or conditions, we sought to

772    calibrate expectations using a large collection of heterogeneous data. This objective

773    aligned with our second aim of identifying reproducible transcription regulator

774    coexpression patterns. We wished to assess how well this information aligns with other

775    lines of regulation evidence, and to provide an organized summary of this information as

776    a community resource (https://doi.org/10.5683/SP3/HJ1B24).

777    While prior work has nominated TR-target interactions across a large and context-

778    independent corpus of data (Garcia-Alonso et al., 2019; Keenan et al., 2019; Müller-Dott

779    et al., 2023), to our knowledge ours is the first to do so using a broad range of single

780    cell transcriptomics. Our literature curation benchmark strongly supports the ability of

781    this resource to prioritize curated targets, and we further find numerous examples of

782    reproducible and conserved coexpressed TR-gene partners also supported by ChIP-

783    seq evidence. Collectively, this suggests that this information can help prioritize

784    interactions when direct experimental evidence is lacking. Our benchmarks additionally

785    provide insight into the TRs whose activity is more challenging to uncover, given the

786    considered genomics data (Supplemental Tables 2-3).

787 Our workflow prioritizes interactions that are most common across contexts, akin to our
788 prior study (Morin et al., 2023). Overall, it is not surprising that the most reproducible
789 relationships tend to relate to processes shared by many cell types. This may be partly
790 a function of expression levels (Crow et al., 2016), but it is logical that the dynamics of
791 processes like the cell cycle are more readily captured by changing transcript levels. We
792 still find evidence for highly context-specific interactions: as long as there is enough
793 supporting data such patterns can emerge. Conversely, if a TR's activity is highly
794 pleiotropic, our framework will tend to only prioritize the partners shared across data.
795 That we are able to observe reproducible patterns in this heterogenous collection raises
796 our confidence in applying this framework to specific contexts in future work, such as
797 identifying tissue-specific versus global partners for TRs like PAX6.

798 Repression is more difficult to infer from coexpression than activation, for reasons we
799 discuss in the Supplemental Material. Similarly, differential interactions are more difficult
800 to characterize than those that are reproducible, requiring evidence of absence. While
801 these considerations motivated our focus on the top reproducible coexpression
802 patterns, the data we have organized can help potentiate the discovery of divergent
803 regulatory interactions. Suresh and colleagues (2023), for example, used single cell
804 coexpression of human and primate data to nominate both conserved and human-novel
805 coexpression patterns. Given that "TR-rewiring" (differential TR activity) is hypothesized
806 to be a primary driver of phenotypic variation, it would be valuable to assess the degree
807 to which differential coexpression between species in matched contexts can reveal
808 distinct regulatory activity.

809 Numerous methods have been developed for gene regulatory network reconstruction
810 using single cell coexpression, with multiple benchmarks concluding that no algorithm
811 dominates (Chen and Mar, 2018; Pratapa et al., 2020; Nguyen et al., 2021; McCalla et
812 al., 2023). In particular, McCalla and colleagues (2023) emphasized the favorable
813 performance of Pearson's correlation (as used in this study) relative to more complex
814 models. This aligns with observations by Harris et al., 2021, who found that aggregating
815 single cell coexpression using the computationally efficient Pearson's correlation
816 provided results that were consistent with alternative similarity metrics (Skinnider et al.,
817 2019). Indeed, we feel that the most important ingredient in the analysis is the
818 aggregation of data because the sparsity of the data is difficult to address otherwise.
819 Our focus on simplistic approaches supports that our conclusions are generalizable to
820 more complex forms of coexpression analysis (Crow et al., 2016).

821 We believe that the organized information we provide will be a valuable community
822 resource. Beyond lists of genes plausibly regulated by each TR, the interactions
823 identified in this study can assist studies examining the conservation of regulatory
824 interactions, or the chromatin factors commonly coexpressed with each TR. Highly
825 ranked interactions could be used for benchmarking predictive methods, or further
826 dissected towards our understanding of the chromatin and sequence features that are
827 characteristic of reproducible interactions. Future work may find it fruitful to construct
828 context-specific aggregations to contrast against this heterogeneous collection, or to
829 further integrate this resource with other lines of regulation evidence, as we did with the
830 ChIP-seq data.

25

## Data Availability

831 All summarized rankings, scored ChIP-seq experiments, and GO analysis results are
832 made available as R objects in the Borealis data repository
833 (https://doi.org/10.5683/SP3/HJ1B24). The identifiers and associated data links of the
834 analyzed scRNA-seq experiments are found in Supplemental Table 1 and summaries of
835 the curation benchmark are found in Supplemental Tables 2-3.  The code to reproduce
836 the analysis is located at https://github.com/PavlidisLab/TR_singlecell.
837

## Funding

## Competing interest

850 The authors have no competing interests to declare.

## Acknowledgments

## Author contributions

860 A.M. conceived and designed the analysis, conducted data collection and analysis, and
861 wrote the manuscript. C.C. contributed to methodological development and data
862 analysis, as well as input on the manuscript. P.P. provided oversight, contributed to
863 study conception and design, and co-wrote the manuscript.
864

## Citations

26

866     1. Aibar S, González-Blas CB, Moerman T, Huynh-Thu VA, Imrichova H, Hulselmans G, et al.
867     SCENIC: single-cell regulatory network inference and clustering. Nature Methods. 2017
868     Nov;14(11):1083–6.
869     2. Al-Jaberi N, Lindsay S, Sarma S, Bayatti N, Clowry GJ. The Early Fetal Development of
870     Human Neocortical GABAergic Interneurons. Cerebral Cortex. 2015 Mar;25(3):631–45.
871     3. Andersson T, Södersten E, Duckworth JK, Cascante A, Fritz N, Sacchetti P, et al. CXXC5 is
872     a novel BMP4-regulated modulator of Wnt signaling in neural stem cells. J Biol Chem. 2009
873     Feb 6;284(6):3672–81.
874     4. Aslanpour S, Rosin JM, Balakrishnan A, Klenin N, Blot F, Gradwohl G, et al. Ascl1 is
875     required to specify a subset of ventromedial hypothalamic neurons. Development. 2020 Jan
876     1;dev.180067.
877     5. Baliñas-Gavira C, Rodríguez MI, Andrades A, Cuadros M, Álvarez-Pérez JC, Álvarez-Prado
878     ÁF, et al. Frequent mutations in the amino-terminal domain of BCL7A impair its tumor
879     suppressor role in DLBCL. Leukemia. 2020 Oct;34(10):2722–35.
880     6. Ballouz S, Verleyen W, Gillis J. Guidance for RNA-seq co-expression network construction
881     and analysis: safety in numbers. Bioinformatics. 2015 Feb 24;btv118.
882     7. Ballouz S, Pavlidis P, Gillis J. Using predictive specificity to determine when gene set
883     analysis is biologically meaningful. Nucleic Acids Res. 2017 Feb 28;45(4):e20–e20.
884     8. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, et al. NCBI GEO:
885     archive for functional genomics data sets--update. Nucleic Acids Res. 2013 Jan;41(Database
886     issue):D991-995.
887     9. Borromeo MD, Savage TK, Kollipara RK, He M, Augustyn A, Osborne JK, et al. ASCL1 and
888     NEUROD1 Reveal Heterogeneity in Pulmonary Neuroendocrine Tumors and Regulate Distinct
889     Genetic Programs. Cell Rep. 2016 Aug 2;16(5):1259–72.
890     10. Bravo González-Blas C, De Winter S, Hulselmans G, Hecker N, Matetovici I, Christiaens V,
891     et al. SCENIC+: single-cell multiomic inference of enhancers and gene regulatory networks.
892     Nat Methods. 2023 Sep;20(9):1355–67.
893     11. Breitling R, Armengaud P, Amtmann A, Herzyk P. Rank products: a simple, yet powerful,
894     new method to detect differentially regulated genes in replicated microarray experiments.
895     FEBS Lett. 2004 Aug 27;573(1–3):83–92.
896     12. Castro DS, Martynoga B, Parras C, Ramesh V, Pacary E, Johnston C, et al. A novel
897     function of the proneural factor Ascl1 in progenitor proliferation identified by genome-wide
898     characterization of its targets. Genes Dev. 2011 May 1;25(9):930–45.
899     13. Castro DS, Skowronska-Krawczyk D, Armant O, Donaldson IJ, Parras C, Hunt C, et al.
900     Proneural bHLH and Brn Proteins Coregulate a Neurogenic Program through Cooperative
901     Binding to a Conserved DNA Motif. Developmental Cell. 2006 Dec 1;11(6):831–44.
902     14. Chen S, Mar JC. Evaluating methods of inferring gene regulatory networks highlights their
903     lack of performance for single cell gene expression data. BMC Bioinformatics. 2018
904     19;19(1):232.
905     15. Chu ECP, Morin A, Chang THC, Nguyen T, Tsai YC, Sharma A, et al. Experiment level
906     curation of transcriptional regulatory interactions in neurodevelopment. PLOS Computational
907     Biology. 2021 Oct 19;17(10):e1009484.
908     16. Crow M, Gillis J. Co-expression in Single-Cell Analysis: Saving Grace or Original Sin?
909     Trends in Genetics [Internet]. 2018 Aug 23 [cited 2018 Aug 28]; Available from:
910     http://www.sciencedirect.com/science/article/pii/S0168952518301288
911     17. Crow M, Paul A, Ballouz S, Huang ZJ, Gillis J. Exploiting single-cell expression to
912     characterize co-expression replicability. Genome Biology. 2016;17:101.
913     18. Cusanovich DA, Pavlovic B, Pritchard JK, Gilad Y. The functional consequences of
914     variation in transcription factor binding. PLoS Genet. 2014 Mar;10(3):e1004226.

915  19. Deprez M, Zaragosi LE, Truchi M, Becavin C, Ruiz García S, Arguel MJ, et al. A Single-Cell
916  Atlas of the Human Healthy Airways. Am J Respir Crit Care Med. 2020 Dec 15;202(12):1636–
917  45.
918  20. Domínguez Conde C, Xu C, Jarvis LB, Rainbow DB, Wells SB, Gomes T, et al. Cross-
919  tissue immune cell analysis reveals tissue-specific features in humans. Science. 2022 May
920  13;376(6594):eabl5197.
921  21. Elmentaite R, Kumasaka N, Roberts K, Fleming A, Dann E, King HW, et al. Cells of the
922  human intestinal tract mapped across space and time. Nature. 2021 Sep 9;597(7875):250–5.
923  22. Emanuele MJ, Enrico TP, Mouery RD, Wasserman D, Nachum S, Tzur A. Complex
924  Cartography: Regulation of E2F Transcription Factors by Cyclin F and Ubiquitin. Trends Cell
925  Biol. 2020 Aug;30(8):640–52.
926  23. Farahbod M, Pavlidis P. Differential coexpression in human tissues and the confounding
927  effect of mean expression levels. Bioinformatics. 2019 Jan 1;35(1):55–61.
928  24. Farahbod M, Pavlidis P. Untangling the effects of cellular composition on coexpression
929  analysis. Genome Res. 2020 Jun 24;30(6):gr.256735.119.
930  25. Fawkner-Corbett D, Antanaviciute A, Parikh K, Jagielowicz M, Gerós AS, Gupta T, et al.
931  Spatiotemporal analysis of human intestinal development at single-cell resolution. Cell. 2021
932  Feb;184(3):810-826.e23.
933  26. Ferreira R, Ohneda K, Yamamoto M, Philipsen S. GATA1 function, a paradigm for
934  transcription factors in hematopoiesis. Mol Cell Biol. 2005 Feb;25(4):1215–27.
935  27. Garcia-Alonso L, Holland CH, Ibrahim MM, Turei D, Saez-Rodriguez J. Benchmark and
936  integration of resources for the estimation of human transcription factor activities. Genome
937  Res. 2019 Aug;29(8):1363–75.
938  28. Garcia-Alonso L, Lorenzi V, Mazzeo CI, Alves-Lopes JP, Roberts K, Sancho-Serra C, et al.
939  Single-cell roadmap of human gonadal development. Nature. 2022 Jul 21;607(7919):540–7.
940  29. Gitter A, Siegfried Z, Klutstein M, Fornes O, Oliva B, Simon I, et al. Backup in gene
941  regulatory networks explains differences between binding and knockout results. Mol Syst Biol.
942  2009;5:276.
943  30. Hamed AA, Kunz DJ, El-Hamamy I, Trinh QM, Subedar OD, Richards LM, et al. A brain
944  precursor atlas reveals the acquisition of developmental-like states in adult cerebral tumours.
945  Nat Commun. 2022 Jul 19;13(1):4178.
946  31. Harris BD, Crow M, Fischer S, Gillis J. Single-cell co-expression analysis reveals that
947  transcriptional modules are shared across cell types in the brain. Cell Systems [Internet]. 2021
948  May 19 [cited 2021 May 25]; Available from:
949  https://www.sciencedirect.com/science/article/pii/S2405471221001538
950  32. Henke RM, Meredith DM, Borromeo MD, Savage TK, Johnson JE. Ascl1 and Neurog2 form
951  novel complexes and regulate Delta-like3 (Dll3) expression in the neural tube. Dev Biol. 2009
952  Apr 15;328(2):529–40.
953  33. Heumos L, Schaar AC, Lance C, Litinetskaya A, Drost F, Zappia L, et al. Best practices for
954  single-cell analysis across modalities. Nat Rev Genet. 2023 Aug;24(8):550–72.
955  34. Hu Y, Comjean A, Rodiger J, Chen W, Kim AR, Qadiri M, et al. FlyRNAi.org 2025 update-
956  expanded resources for new technologies and species. Nucleic Acids Res. 2025 Jan
957  6;53(D1):D958–65.
958  35. Hu Y, Flockhart I, Vinayagam A, Bergwitz C, Berger B, Perrimon N, et al. An integrative
959  approach to ortholog prediction for disease-focused and other functional studies. BMC
960  Bioinformatics. 2011 Aug 31;12:357.
961  36. Hu Z, Killion PJ, Iyer VR. Genetic reconstruction of a functional transcriptional regulatory
962  network. Nat Genet. 2007 May;39(5):683–7.
963  37. Huang HS, Kubish GM, Redmond TM, Turner DL, Thompson RC, Murphy GG, et al. Direct
964  transcriptional induction of Gadd45gamma by Ascl1 during neuronal differentiation. Mol Cell
965  Neurosci. 2010 Jul;44(3):282–96.

966 38. Jacob J, Storm R, Castro DS, Milton C, Pla P, Guillemot F, et al. Insm1 (IA-1) is an
967 essential component of the regulatory network that specifies monoaminergic neuronal
968 phenotypes in the vertebrate hindbrain. Development. 2009 Jul;136(14):2477–85.
969 39. Jia S, Wildner H, Birchmeier C. Insm1 controls the differentiation of pulmonary
970 neuroendocrine cells by repressing Hes1. Dev Biol. 2015 Dec 1;408(1):90–8.
971 40. Kang Y, Patel NR, Shively C, Recio PS, Chen X, Wranik BJ, et al. Dual threshold
972 optimization and network inference reveal convergent evidence from TF binding locations and
973 TF perturbation responses. Genome Res. 2020 Mar;30(3):459–71.
974 41. Keenan AB, Torre D, Lachmann A, Leong AK, Wojciechowicz ML, Utti V, et al. ChEA3:
975 transcription factor enrichment analysis by orthogonal omics integration. Nucleic Acids Res.
976 2019 Jul 2;47(W1):W212–24.
977 42. Kito-Shingaki A, Seta Y, Toyono T, Kataoka S, Kakinoki Y, Yanagawa Y, et al. Expression
978 of GAD67 and Dlx5 in the Taste Buds of Mice Genetically Lacking Mash1. Chemical Senses.
979 2014 Jun 1;39(5):403–14.
980 43. Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, et al. The Human
981 Transcription Factors. Cell. 2018 Feb 8;172(4):650–65.
982 44. Lampada A, Taylor V. Notch signaling as a master regulator of adult neurogenesis. Front
983 Neurosci. 2023 Jun 29;17:1179011.
984 45. Lee HK, Hsu AK, Sajdak J, Qin J, Pavlidis P. Coexpression analysis of human genes
985 across many microarray data sets. Genome Res. 2004;14:1085–94.
986 46. Lee J, Shah M, Ballouz S, Crow M, Gillis J. CoCoCoNet: conserved and comparative co-
987 expression across a diverse set of species. Nucleic Acids Res. 2020 Jul 2;48(W1):W566–71.
988 47. Li X, Zheng Y, Hu H, Li X. Integrative analyses shed new light on human ribosomal protein
989 gene regulation. Sci Rep [Internet]. 2016 Jun 27 [cited 2018 Dec 25];6. Available from:
990 https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4921865/
991 48. Lin H, Zhu X, Chen G, Song L, Gao L, Khand AA, et al. KDM3A-mediated demethylation of
992 histone H3 lysine 9 facilitates the chromatin binding of Neurog2 during neurogenesis.
993 Development. 2017 Oct 15;144(20):3674–85.
994 49. Liu C, Martins AJ, Lau WW, Rachmaninoff N, Chen J, Imberti L, et al. Time-resolved
995 systems immunology reveals a late juncture linked to fatal COVID-19. Cell. 2021 Apr
996 1;184(7):1836-1857.e22.
997 50. Liu J, Gao L, Ji B, Geng R, Chen J, Tao X, et al. BCL7A as a novel prognostic biomarker
998 for glioma patients. J Transl Med. 2021 Aug 6;19(1):335.
999 51. Liu YH, Tsai JW, Chen JL, Yang WS, Chang PC, Cheng PL, et al. Ascl1 promotes
1000 tangential migration and confines migratory routes by induction of Ephb2 in the telencephalon.
1001 Sci Rep. 2017 Mar 9;7(1):42895.
1002 52. Lunden JW, Durens M, Phillips AW, Nestor MW. Cortical interneuron function in autism
1003 spectrum condition. Pediatr Res. 2019 Jan;85(2):146–54.
1004 53. Marbach D, Costello JC, Küffner R, Vega NM, Prill RJ, Camacho DM, et al. Wisdom of
1005 crowds for robust gene network inference. Nature Methods [Internet]. 2012 [cited 2012 Jul 18];
1006 Available from: http://www.nature.com/nmeth/journal/vaop/ncurrent/abs/nmeth.2016.html
1007 54. Mathys H, Davila-Velderrain J, Peng Z, Gao F, Mohammadi S, Young JZ, et al. Single-cell
1008 transcriptomic analysis of Alzheimer's disease. Nature. 2019 Jun 20;570(7761):332–7.
1009 55. Mazurier N, Parain K, Parlier D, Pretto S, Hamdache J, Vernier P, et al. Ascl1 as a novel
1010 player in the Ptf1a transcriptional network for GABAergic cell specification in the retina. PLoS
1011 One. 2014;9(3):e92113.
1012 56. McCall MN, Illei PB, Halushka MK. Complex Sources of Variation in Tissue Expression
1013 Data: Analysis of the GTEx Lung Transcriptome. The American Journal of Human Genetics.
1014 2016 Sep 1;99(3):624–35.
1015 57. McCalla SG, Fotuhi Siahpirani A, Li J, Pyne S, Stone M, Periyasamy V, et al. Identifying
1016 strengths and weaknesses of methods for computational network inference from single-cell

1017    RNA-seq data. Steinmetz L, editor. G3: Genes, Genomes, Genetics. 2023 Mar
1018    9;13(3):jkad004.
1019    58. Mistry M, Gillis J, Pavlidis P. Meta-analysis of gene coexpression networks in the post-
1020    mortem prefrontal cortex of patients with schizophrenia and unaffected controls. BMC
1021    Neurosci. 2013 Dec;14(1):105.
1022    59. Monaco G, Van Dam S, Casal Novo Ribeiro JL, Larbi A, De Magalhães JP. A comparison
1023    of human and mouse gene co-expression networks reveals conservation and divergence at the
1024    tissue, pathway and disease levels. BMC Evol Biol. 2015 Dec;15(1):259.
1025    60. Morin A, Chu ECP, Sharma A, Adrian-Hamazaki A, Pavlidis P. Characterizing the targets of
1026    transcription regulators by aggregating ChIP-seq and perturbation expression data sets.
1027    Genome Res [Internet]. 2023 Jun 12 [cited 2023 Jun 12]; Available from:
1028    https://genome.cshlp.org/content/early/2023/06/12/gr.277273.122
1029    61. Müller-Dott S, Tsirvouli E, Vazquez M, Ramirez Flores RO, Badia-I-Mompel P, Fallegger R,
1030    et al. Expanding the coverage of regulons from high-confidence prior knowledge for accurate
1031    estimation of transcription factor activities. Nucleic Acids Res. 2023 Nov 10;51(20):10934–49.
1032    62. Nelson BR, Hartman BH, Ray CA, Hayashi T, Bermingham-McDonogh O, Reh TA.
1033    Acheate-scute like 1 (Ascl1) is required for normal delta-like (Dll) gene expression and notch
1034    signaling during retinal development. Dev Dyn. 2009 Sep;238(9):2163–78.
1035    63. Nguyen H, Tran D, Tran B, Pehlivan B, Nguyen T. A comprehensive survey of regulatory
1036    network inference methods using single cell RNA sequencing data. Briefings in Bioinformatics.
1037    2021 May 20;22(3):bbaa190.
1038    64. Ouyang Z, Zhou Q, Wong WH. ChIP-Seq of transcription factors predicts absolute and
1039    differential gene expression in embryonic stem cells. Proc Natl Acad Sci USA. 2009 Dec
1040    22;106(51):21521–6.
1041    65. Patel RV, Nahal HK, Breit R, Provart NJ. BAR expressolog identification: expression profile
1042    similarity ranking of homologous genes in plant species. Plant J. 2012 Sep;71(6):1038–50.
1043    66. Posner DA, Lee CY, Portet A, Clatworthy MR. Humoral immunity at the brain borders in
1044    homeostasis. Current Opinion in Immunology. 2022 Jun;76:102188.
1045    67. Pratapa A, Jalihal AP, Law JN, Bharadwaj A, Murali TM. Benchmarking algorithms for gene
1046    regulatory network inference from single-cell transcriptomic data. Nat Methods.
1047    2020;17(2):147–54.
1048    68. Puig RR, Boddie P, Khan A, Castro-Mondragon JA, Mathelier A. UniBind: maps of high-
1049    confidence direct TF-DNA interactions across nine species. BMC Genomics. 2021
1050    Dec;22(1):482.
1051    69. Qu S, Fetsch P, Thomas A, Pommier Y, Schrump DS, Miettinen MM, et al. Molecular
1052    Subtypes of Primary SCLC Tumors and Their Associations With Neuroendocrine and
1053    Therapeutic Markers. J Thorac Oncol. 2022 Jan;17(1):141–53.
1054    70. Ragazzini R, Boeing S, Zanieri L, Green M, D'Agostino G, Bartolovic K, et al. Defining the
1055    identity and the niches of epithelial stem cells with highly pleiotropic multilineage potency in the
1056    human thymus. Developmental Cell. 2023 Nov;58(22):2428-2446.e9.
1057    71. Rothenberg EV. Causal Gene Regulatory Network Modeling and Genomics: Second-
1058    Generation Challenges. Journal of Computational Biology. 2019 Jul 1;26(7):703–18.
1059    72. Russo GL, Sonsalla G, Natarajan P, Breunig CT, Bulli G, Merl-Pham J, et al. CRISPR-
1060    Mediated Induction of Neuron-Enriched Mitochondrial Proteins Boosts Direct Glia-to-Neuron
1061    Conversion. Cell Stem Cell. 2021 Mar 4;28(3):524-534.e7.
1062    73. Salem S, Salem D, Gros P. Role of IRF8 in immune cells functions, protection against
1063    infections, and susceptibility to inflammatory diseases. Hum Genet. 2020 Jun;139(6–7):707–
1064    21.
1065    74. Shen WK, Chen SY, Gan ZQ, Zhang YZ, Yue T, Chen MM, et al. AnimalTFDB 4.0: a
1066    comprehensive animal transcription factor database updated with variation and expression
1067    annotations. Nucleic Acids Research. 2023 Jan 6;51(D1):D39–45.

75. Simão D, Silva MM, Terrasso AP, Arez F, Sousa MFQ, Mehrjardi NZ, et al. Recapitulation of Human Neural Microenvironment Signatures in iPSC-Derived NPC 3D Differentiation. Stem Cell Reports. 2018 Aug;11(2):552–64.

76. Sing T, Sander O, Beerenwinkel N, Lengauer T. ROCR: visualizing classifier performance in R. Bioinformatics. 2005 Oct 15;21(20):3940–1.

77. Skinnider MA, Squair JW, Foster LJ. Evaluating measures of association for single-cell transcriptomics. Nature Methods. 2019 Apr 8;1.

78. Sonawane AR, Weiss ST, Glass K, Sharma A. Network Medicine in the Age of Biomedical Big Data. Front Genet. 2019;10:294.

79. Suresh H, Crow M, Jorstad N, Hodge R, Lein E, Dobin A, et al. Comparative single-cell transcriptomic analysis of primate brains highlights human-specific regulatory evolution. Nat Ecol Evol. 2023 Sep 4;1–14.

80. Szu J, Wojcinski A, Jiang P, Kesari S. Impact of the Olig Family on Neurodevelopmental Disorders. Front Neurosci. 2021 Mar 30;15:659601.

81. Tamrazi B, Venneti S, Margol A, Hawes D, Cen SY, Nelson M, et al. Pediatric Atypical Teratoid/Rhabdoid Tumors of the Brain: Identification of Metabolic Subgroups Using In Vivo 1H-MR Spectroscopy. AJNR Am J Neuroradiol. 2019 May;40(5):872–7.

82. Tutukova S, Tarabykin V, Hernandez-Miranda LR. The Role of Neurod Genes in Brain Development, Function, and Disease. Front Mol Neurosci. 2021 Jun 9;14:662774.

83. Ueno T, Ito J, Hoshikawa S, Ohori Y, Fujiwara S, Yamamoto S, et al. The identification of transcriptional targets of Ascl1 in oligodendrocyte development. Glia. 2012 Oct;60(10):1495–505.

84. Uhlén M, Fagerberg L, Hallström BM, Lindskog C, Oksvold P, Mardinoglu A, et al. Tissue-based map of the human proteome. Science. 2015 Jan 23;347(6220):1260419.

85. van Lengerich B, Zhan L, Xia D, Chan D, Joy D, Park JI, et al. A TREM2-activating antibody with a blood-brain barrier transport vehicle enhances microglial metabolism in Alzheimer's disease models. Nat Neurosci. 2023 Mar;26(3):416–29.

86. Wang B, Long JE, Flandin P, Pla R, Waclaw RR, Campbell K, et al. Loss of Gsx1 and Gsx2 function rescues distinct phenotypes in Dlx1/2 mutants. J of Comparative Neurology. 2013 May;521(7):1561–84.

87. Wang S, Sun H, Ma J, Zang C, Wang C, Wang J, et al. Target analysis by integration of transcriptome and ChIP-seq data with BETA. Nat Protoc. 2013 Dec;8(12):2502–15.

88. Wang X, He Y, Zhang Q, Ren X, Zhang Z. Direct Comparative Analyses of 10X Genomics Chromium and Smart-seq2. Genomics, Proteomics & Bioinformatics. 2021 Apr;19(2):253–66.

89. Wen JH, Chen YY, Song SJ, Ding J, Gao Y, Hu QK, et al. Paired box 6 (PAX6) regulates glucose metabolism via proinsulin processing mediated by prohormone convertase 1/3 (PC1/3). Diabetologia. 2009 Mar 1;52(3):504–13.

90. Weng Q, Wang J, Wang J, He D, Cheng Z, Zhang F, et al. Single-Cell Transcriptomics Uncovers Glial Progenitor Diversity and Cell Fate Determinants during Development and Gliomagenesis. Cell Stem Cell. 2019 May;24(5):707-723.e8.

91. Werner JM, Gillis J. Preservation of co-expression defines the primary tissue fidelity of human neural organoids. bioRxiv. 2023 Oct 17;2023.03.31.535112.

92. Xiong X, Tu S, Wang J, Luo S, Yan X. CXXC5: A novel regulator and coordinator of TGF-β, BMP and Wnt signaling. J Cell Mol Med. 2019 Feb;23(2):740–9.

93. Yamada Y, Bohnenberger H, Kriegsmann M, Kriegsmann K, Sinn P, Goto N, et al. Tuft cell-like carcinomas: novel cancer subsets present in multiple organs sharing a unique gene expression signature. Br J Cancer. 2022 Nov;127(10):1876–85.

94. Yeung J, Ha TJ, Swanson DJ, Goldowitz D. A Novel and Multivalent Role of Pax6 in Cerebellar Development. J Neurosci. 2016 Aug 31;36(35):9057–69.

95. Yu J, Mu J, Guo Q, Yang L, Zhang J, Liu Z, et al. Transcriptomic profile analysis of mouse neural tube development by RNA-Seq. IUBMB Life. 2017 Jul 10;

1119    96. Yue F, Cheng Y, Breschi A, Vierstra J, Wu W, Ryba T, et al. A comparative encyclopedia of
1120    DNA elements in the mouse genome. Nature. 2014 Nov 20;515(7527):355–64.
1121    97. Zhang W, Girard L, Zhang YA, Haruki T, Papari-Zareei M, Stastny V, et al. Small cell lung
1122    cancer tumors and preclinical models display heterogeneity of neuroendocrine phenotypes.
1123    Transl Lung Cancer Res. 2018 Feb;7(1):32–49.
1124    98. Zhang Y, Cuerdo J, Halushka MK, McCall MN. The effect of tissue composition on gene
1125    co-expression. Brief Bioinform [Internet]. [cited 2019 Dec 9]; Available from:
1126    https://academic.oup.com/bib/advance-article/doi/10.1093/bib/bbz135/5669861
1127    99. Zhong E, Pareja F, Hanna MG, Jungbluth AA, Rekhtman N, Brogi E. Expression of novel
1128    neuroendocrine markers in breast carcinomas: a study of INSM1, ASCL1, and POU2F3. Hum
1129    Pathol. 2022 Sep;127:102–11.
1130    100. Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. Nature. 2018 Oct
1131    3;1.