

1 *Acorde*: unraveling functionally-interpretable networks of isoform co-usage
2 from single cell data

3 Angeles Arzalluz-Luque¹, Pedro Salguero¹, Sonia Tarazona^{1,+,*}, Ana Conesa^{2,+,*}

4 ¹ Department of Applied Statistics, Operations Research and Quality, Universitat Politècnica de València, Valencia, Spain.

5 ² Microbiology and Cell Sciences Department, Institute for Food and Agricultural Research, University of Florida,
6 Gainesville, Florida, United States.

7 ³ Genetics Institute, University of Florida, Gainesville, Florida, United States.

8 ⁺ These authors jointly supervised this work.

9 ^{*} Corresponding authors: sotacam@eio.upv, aconesa@ufl.edu

10 **Abstract**

11 Alternative splicing (AS) is a highly-regulated post-transcriptional mechanism known to modulate isoform
12 expression within genes and contribute to cell-type identity. However, the extent to which alternative isoforms
13 establish co-expression networks that may relevant in cellular function has not been explored yet. Here, we
14 present *acorde*, a pipeline that successfully leverages bulk long reads and single-cell data to confidently detect
15 alternative isoform co-expression relationships. To achieve this, we developed and validated percentile
16 correlations, a novel approach that overcomes data sparsity and yields accurate co-expression estimates from
17 single-cell data. Next, *acorde* uses correlations to cluster co-expressed isoforms into a network, unraveling cell
18 type-specific alternative isoform usage patterns. By selecting same-gene isoforms between these clusters, we
19 subsequently detect and characterize genes with co-differential isoform usage (coDIU) across neural cell types.
20 Finally, we predict functional elements from long read-defined isoforms and provide insight into biological
21 processes, motifs and domains potentially controlled by the coordination of post-transcriptional regulation.

1 **Introduction**

2 Single-cell RNA-seq (scRNA-seq) has revolutionized transcriptomics analysis, especially as the
3 development of technologies with increasingly higher throughputs has enabled the
4 processing of thousands of single cells simultaneously, boosting the amount of biological
5 diversity that can be captured in a sequencing experiment¹. The technology has been
6 extensively applied to the discovery of new cell types and the characterization of their
7 transcriptional profiles, resulting in the definition of cell type marker genes²⁻⁷. Even though
8 scRNA-seq datasets can encode several levels of granularity in the form of cell subtypes, these
9 studies rely on the low number of features required to recapitulate the cell type structure of
10 the data⁸, which situates cell type characterization efforts at the baseline of understanding the
11 intricacy of single-cell biology. scRNA-seq studies have also tackled transcriptional dynamics
12 and how they relate to cell type properties. These include methods for the study of dynamic
13 processes, namely pseudotime^{9,10} and RNA-velocity¹¹ analyses, which have provided insight on
14 cell differentiation and transition mechanisms between cell states^{7,12-16}. Moreover, novel
15 methods have been developed to convey the inference of gene regulatory networks (GRNs) to
16 the single-cell level¹⁷ in an attempt to combine single-cell information with extant knowledge
17 to infer relationships between genes and transcriptional regulators at a higher resolution.

18 Single-cell research is nevertheless far from realizing its full potential. On the contrary, the
19 timing is now optimal for the field to undertake the investigation of deeper layers of cellular
20 complexity. In particular, the investigation of Alternative Splicing (AS) and isoform expression
21 dynamics has remained a challenge to the field. Reasons for this include the uncertainty of
22 short read-based isoform quantification, which is exacerbated by the lower number of
23 available reads per transcript in comparison to bulk libraries¹⁸, and the fact that the most

1 popular scRNA-seq methods are heavily 3' end biased, which precludes the unambiguous
2 identification of alternative transcript variants¹⁹. Current methods for the study of isoforms in
3 single cells therefore rely on alternative metrics that either avoid isoform-level expression
4 estimation completely^{20,21} or exclusively consider individual splicing events²²⁻²⁶, generally
5 leaving isoform characterization aside, with some recent exceptions²⁷. Meanwhile, long read
6 RNA sequencing (lrRNA-seq) of single cells is beginning to emerge an alternative approach to
7 mitigate this ambiguity, given that it successfully grasps how individual events are combined
8 into alternative isoforms¹⁹. Long read studies have therefore expanded the field's notions of
9 cell type-specific splicing from event inclusion towards isoform selection patterns and showed
10 that cell type-specific isoform expression can be detected in both broad types as well as cell
11 subtypes²⁸⁻³². Unfortunately, the sequencing depth constraints intrinsic to long read
12 protocols¹⁹ have limited the amount of isoform diversity that can be captured by single-cell
13 long read transcriptomics²⁸⁻³⁰, and datasets generally show low levels of redundancy between
14 cells of the same cell type.

15 Notwithstanding this limited scenario, there are a number of relevant questions regarding the
16 importance of splicing for cell identity and function that can only be resolved by evaluating
17 isoform expression at the single-cell level. In point of fact, splicing differences have been
18 shown to discriminate cell types with an accuracy comparable to that obtained using gene
19 expression³³, while integrating AS and gene expression changes has led to the discovery of cell
20 subtypes and states that were otherwise not detected^{27,34-36}. Especially relevant among these
21 inquiries is the much-debated issue of whether individual cells express one or several
22 isoforms, that is, whether the isoform diversity observed in bulk studies is recapitulated by
23 each single cell or, alternatively, arises as a result of the combination of multiple cells, each of

1 which uniquely express one of the gene's isoforms. Ever since the publication of the first
2 scRNA-seq studies, short reads have been used to answer this question, usually via the
3 characterization of splicing event -rather than isoform expression- modalities. Successive
4 studies have provided non-conclusive results, with evidence of bimodal splicing patterns^{22,37,38}
5 as well as concerns regarding the relationship between bimodal isoform detection and
6 technical noise^{18,39}, a controversy that suggests that new analytical and computational
7 approaches are needed to understand the isoform landscape of single cells.

8 Another pending question for the field is whether isoform expression programs involve co-
9 expression relationships between transcript variants from different genes. So far, the
10 application of long read technologies to single-cell data has served to unravel coordinated
11 event choice patterns within isoforms of the same gene^{40,41}, however, cross-gene isoform
12 expression networks have not been investigated. In other words, there have been no studies
13 addressing potential codependency between genes regarding the selection of transcript
14 variants from their isoform repertoire, or the implications of this coordination for cell-state
15 and cell-type properties. This is not only related to the general constraints of single-cell
16 isoform studies, but also to the lack of computational methods and mathematical models to
17 extract this complex signal from the data. In spite of the present research gap, isoform co-
18 expression networks are an anticipated consequence of the regulation of splicing by RNA
19 binding proteins (RBPs) and other splicing factors, and their investigation constitutes an
20 opportunity to gain insight on the functional role of AS. Moreover, given that a multiple cell
21 type and high cell throughput context such as that of scRNA-seq data constitutes a far more
22 suitable data scenario to that of bulk RNA-seq, this is undoubtedly a timely inquiry to make.

1 In the present study, we hypothesize that isoform expression coordination exists as a result of
2 AS regulation, and that this can be computationally detected in the form of isoform groups
3 showing co-variation across cell types. To demonstrate this hypothesis, we have designed an
4 end-to-end, data-intensive pipeline for the study of isoform networks ([Figure 1](#)). First and
5 foremost, we employed a hybrid strategy where bulk long-reads and single-cell Illumina
6 sequencing were integrated to estimate isoform expression at the single cell level. To unlock
7 the limitations of extant correlation metrics in the single-cell context⁴², we developed a novel
8 strategy to obtain noise-robust correlation estimates in scRNA-seq data, and a semi-
9 automated clustering approach to detect modules of co-expressed isoforms across cell types.
10 We additionally re-defined and implemented Differential Isoform Usage (DIU) and co-
11 Differential Isoform Usage (coDIU) analyses in order to leverage the multiple cell types
12 contained in single-cell datasets. Finally, to couple these analyses with a biologically
13 interpretable readout, we incorporated a functional annotation step in which several
14 databases and prediction tools were integrated to add isoform-specific functional information
15 ([Figure 1](#)).

16 We have hereby applied this pipeline ([Figure 1](#)) to the analysis of a publicly available mouse
17 neural dataset, including published scRNA-seq Smart-Seq2 (primary visual cortex, generated
18 by Tasic *et al.*⁴³) and bulk ENCODE PacBio long-read data (mouse cortex and hippocampus,
19 generated by Wyman *et al.*⁴⁴). As a result, we successfully detected cell type-specific co-
20 expression of isoforms in a manner that was independent of gene-level expression.
21 Furthermore, we demonstrated that these isoforms encode shared functional properties,
22 highlighting the role of post-transcriptional processing as an additional regulatory layer that

- 1 fine-tunes cellular functions and contributes to encode cell type identity. This novel pipeline
- 2 has been implemented in the R package *acorde* (<https://github.com/ConesaLab/acorde>).

3 **Results**

4 *Enabling multi-group differential expression of isoforms in single-cell data*

5 Traditionally, RNA-seq studies use publicly available reference transcriptomes such as RefSeq
6 and ENSEMBL for short read isoform quantification. However, most tissues and cell types will
7 express only a subset of the genes and isoforms contained in the reference, with previous
8 studies showing that isoform detection accuracy increases when adopting tissue-specific
9 isoform sets as a reference for mapping⁴⁵. Long read technologies have the potential to
10 achieve this, while also expanding the reference with novel isoforms that have not yet been
11 annotated⁴⁶. Given the low depth of single-cell long read datasets¹⁹, we employed a bulk
12 PacBio dataset obtained from ENCODE⁴⁴ to define a mouse, neural-specific transcriptome,
13 that, after extensive curation using the SQANTI3 toolkit
14 (<https://github.com/ConesaLab/SQANTI3>) contained 36,986 isoforms belonging to 12,692
15 genes (see [Supplementary Note 1](#)).

16 To quantify the expression of the long read-defined isoforms at the single-cell level, we made
17 use of a publicly available, deeply sequenced, full-length, short-read single-cell RNA-seq
18 dataset by Tasic *et al.*⁴³. We retained 1,591 cells after quality control (see Methods). Using the
19 labels from the original characterization of the dataset, we assigned cells to 7 broad cell types,
20 5 glial (microglia, endothelial cells, oligodendrocytes, oligodendrocyte precursor cells (OPCs)

1 and astrocytes) and 2 neural (GABA-ergic and glutamatergic neurons), each of which can be
2 divided into several, distinct subtypes ([Supplementary Figure 1A](#)). To target potential
3 differential isoform selection between cell types, we selected genes that retained more than
4 one isoform after quality filtering, ultimately keeping 13,832 isoforms from 4,591 genes.

5 Of note, we observed a drastic cell number imbalance between neural (~720 cells/cell type)
6 and glial cell types (~30 cells/cell type) in the Tasic dataset, which resulted in the
7 underestimation of transcriptional differences between non-neural types ([Supplementary](#)
8 [Figure 1B](#)). In order to facilitate downstream analysis, we used a downsampling approach
9 (Methods) to balance the cell type abundances in the dataset while effectively preserving the
10 data structure ([Supplementary Figure 1C](#)). Next, to select isoforms with robust co-variation
11 and non-constitutive expression, we applied a multi-group strategy to detect isoforms
12 showing Differential Expression (DE) in at least one cell type, combining the ZinBWave zero-
13 expression weighting strategy⁴⁷ with bulk-designed DE methods DESeq2⁴⁸ and edgeR⁴⁹ (see
14 Methods). Upon testing, we detected 5,711 DE isoforms using edgeR (FDR<0.05) and 7,714
15 using DESeq2 (FDR<0.05). To maximize sensitivity and ensure a broad iso-transcriptome
16 analysis, we considered an isoform to be DE if it was detected by at least one of these
17 methods, resulting in 10,100 isoforms from 4,305 genes ([Supplementary Figure 2](#)). Since two
18 or more isoforms with differential cell type expression are required to form co-splicing
19 relationships (see Methods), we finally retained 8,967 isoforms from 3,172 genes with multiple
20 DE isoforms for downstream analysis.

21 To ensure that downsampling had no effect on DE results, a validation experiment for multi-
22 group DE analysis was conducted, in which we performed 50 runs of random neural cell

1 sampling (see Methods). Although DESeq2 proved to be slightly more robust across
2 independent runs than edgeR, Jaccard Index values indicated that the majority of isoforms
3 were consistently detected as DE (DESeq2: mean no. of DE isoforms = $6,908 \pm 101$, mean
4 Jaccard Index = 0.84 ± 0.02 ; edgeR: mean no. of DE isoforms = $6,016 \pm 410$, mean Jaccard Index =
5 0.74 ± 0.02). In addition, this validation run showed that considering the union of edgeR and
6 DESeq2 results contributed to improve robustness (mean no. of DE isoforms in union
7 between methods = $9,399 \pm 248$, Jaccard Index = 0.82 ± 0.01).

8 *Detecting isoforms showing cell type-level co-expression*

9 Co-expression signals in single-cell data are weak and often result in poor performance of
10 traditional correlation metrics and network inference methods^{42,50}. Although data
11 transformation approaches⁵¹ and alternative metrics⁴² have been proposed, these are more
12 complex to apply and considerably less interpretable, respectively. Furthermore, most of
13 these studies have only investigated gene level co-expression¹⁷, often ignoring the AS
14 regulatory landscape. To address these limitations, we implemented a *percentile correlation*
15 strategy: a simple, scalable approach to overcome single-cell noise in isoform co-expression
16 studies (Figure 2A).

17 Our approach considers cell-type identity to be defined by context-specific regulation of gene
18 expression programs, and within-cell type stochasticity to arise from a combination of
19 technical noise^{52,53} and biological mechanisms such as transcriptional bursting⁵⁴, which
20 translate into the sparsity and heterogeneous expression patterns typically observed in
21 scRNA-seq (Supplementary Figure 3A) and result in high variance across all levels of

1 expression ([Supplementary Figure 3B](#)). Together, these effects mask the co-expression signal
2 in the data and tend to yield low correlation values when using traditional metrics ([Figure 2B](#)).
3 To overcome this problem, we treat single cells of the same cell-type as biological replicates
4 that represent the state of a delimited cell population but are differently affected by the
5 aforementioned combination of technical and biological forces. In this context, the expression
6 distribution of any given transcript across the population can be considered as the signature
7 of the transcript in that cell-type. To translate these assumptions into a metric, every isoform's
8 expression within a cell type is first summarized into an expression profile, where single-cell
9 count values are replaced by 10 percentile values (deciles) ([Figure 2A](#), Methods). Intuitively,
10 this reduced number of values synthesizes the approximate behavior of that transcript in the
11 cell type, as inferred given cell-level observations. Next, to grasp expression distribution
12 similarities across cell types, Pearson correlation between transcript pairs is computed using
13 the percentile-summarized expression, resulting in a meaningful distribution of correlation
14 values ([Figure 2B](#)). In this manner, we extend the notion of cell-type markers to rely not only
15 on mean or frequency of expression, but on their actual distributional pattern. Our co-
16 expression metric therefore by-passes cell-level matching of individual observations,
17 providing a correlation estimate that is both robust to the uncertainty of single cell expression
18 and interpretable as a measure of expression similarity. Of note, changing percentile number
19 did not have a noticeable effect on the resulting correlation values ([Supplementary Figure 3C](#)),
20 which were considerably higher than those generated using traditional metrics, successfully
21 unlocking co-expression analysis in single-cell data.

22 To detect modules of co-expressed isoforms, we used the percentile correlation matrix as a
23 distance matrix for hierarchical clustering, and designed a semi-automated cluster refinement

1 approach to ensure maximal profile similarity within clustered modules ([Figure 3A](#), see
2 Methods for a detailed description). Briefly, we first used the *dynamicTreeCut* R package⁵⁵ to
3 set adaptive thresholds and find clusters dynamically within the dendrogram. Among the
4 resulting 152 clusters, some showed highly similar (i.e. redundant) expression profiles
5 ([Supplementary Figure 4A](#)) and were therefore merged by re-clustering using the mean scaled
6 expression of the isoforms in each cluster (i.e. clustering of metatranscripts, see Methods),
7 obtaining 60 isoform clusters. Among them, 18 groups presented noisy expression profiles
8 ([Supplementary Figure 4B](#)) and their isoforms were therefore re-assigned to the remaining 42
9 clusters. However, fully automated re-clustering was not effective to eliminate redundancy in
10 some cases. As a result, conflicting cases were inspected and merged, generating a final
11 number of 19 distinct clusters containing 8,688 isoforms (96,8% of total) and representing
12 diverse expression modalities across the 7 broad cell types ([Figure 3B](#)).

13 *Validation of percentile correlations on simulated data*

14 Building on studies reporting the poor performance of popular correlation metrics in single-
15 cell data, authors have attempted the implementation of sparsity-aware measurements^{51,56}
16 and reported the potential of other alternatives to compute similarity, such as proportionality
17 metrics⁴². Here, we present an interpretable, scalable and biology-aware alternative to single-
18 cell co-expression studies based on Pearson or Spearman correlation. However, to better
19 understand how percentile correlation performs in comparison to extant correlation metrics,
20 we compared it to Pearson, Spearman and zero-inflated Kendall correlations⁵⁶ and one
21 proportionality metric, rho (ρ)⁵⁷ using simulated data.

1 Given that there are -to the best of our knowledge- no scRNA-seq data simulators that include
2 transcript co-expression patterns, we designed a simulation strategy ([Supplementary Figure](#)
3 [5A](#), Methods) to generate an appropriate validation framework for our metric. Briefly, we
4 applied *SymSim*⁵⁸ to simulate a single-cell RNA-seq dataset (8 cell types, 1,000 cells and 8,000
5 transcripts) and used the simulated expression values to artificially create 3,000 synthetic
6 transcripts showing 15 different expression profiles across the 8 cell types ([Supplementary](#)
7 [Figure 5B](#)). As a result, our simulated dataset contained 15 simulated clusters with distinct
8 expression profiles. Among them, clusters 1-5, 6-10 and 11-15 included transcripts showing
9 high expression in one, two and three cell types respectively, gradually increasing simulated
10 pattern complexity ([Supplementary Figures 5B-C](#)). After refinement ([Supplementary Figure 5C](#),
11 Methods), 1,790 synthetic transcripts remained distributed across the 15 simulated clusters in
12 groups ranging from 180 to 60 transcripts ([Supplementary Figure 5D](#)).

13 In order to evaluate how well the 5 co-expression methods recapitulated the simulated
14 patterns, we computed these metrics for all synthetic transcript pairs in each simulated
15 cluster ([Supplementary Figure 5E](#)). Among them, percentile correlation consistently yielded
16 the best proportion of high within-cluster correlations followed by ρ , however, rather counter-
17 intuitively, ρ had only an average performance when low-complexity patterns were provided,
18 with less than 20% output proportionality values >0.8 within clusters 1-5. Shockingly, zero-
19 inflated Kendall correlation, a single cell-tailored metric, failed to recapitulate the simulated
20 co-expression profiles and showed a considerably lower proportion of high correlations
21 within the simulated clusters than Pearson and Spearman correlations. As a result of this
22 evaluation, we can confidently assume that percentile correlations are useful to detect co-
23 variation patterns, yielding overall higher correlation values than all other considered metrics.

1 Of note, and similarly to real data, the simulated dataset showed no detectable effect when
2 varying the number of percentiles used to compute percentile correlations ([Supplementary](#)
3 [Figure 6](#)).

4 Next, we compared the ability of each co-expression metric to inform clustering and group
5 synthetic transcripts with similar expression patterns. To achieve this, we run our clustering
6 pipeline on the simulated isoforms using the 5 metrics as distance ([Supplementary Figures 7-](#)
7 [11](#)). To enable benchmarking, clustering was automated to always generate a total number of
8 15 calculated clusters (see Methods). In order to evaluate which metric worked best to detect
9 co-expressed transcript groups, we considered internal correlations between the transcripts
10 in generated clusters. We observed that ρ and percentile-generated clusters, unlike the
11 remaining co-expression metrics, presented consistently high levels of internal correlation
12 ([Figure 4A](#)). Notably, the distribution of co-expression values obtained using percentiles was
13 the most robust among the five metrics ([Supplementary Figure 12](#)). We next assessed how
14 well the clusters generated using each co-expression metric (i.e. calculated clusters)
15 recapitulated the simulated clusters. Calculated and simulated clusters were paired based on
16 the similarities between their mean cluster profiles ([Supplementary Figure 13](#), Methods), and
17 the Jaccard index (JI) for each simulated-calculated pair was computed to measure the
18 agreement in synthetic transcript assignment ([Figure 4B](#)). Interestingly, results were highly
19 heterogeneous for most methods: even though a number of simulated co-expression groups
20 were easily detected by most metrics, no method was able to fully recapitulate the simulated
21 clusters, with ρ proportionality, Pearson and percentile correlations being the most accurate
22 ([Figure 4B](#)). Zero-inflated Kendall and Spearman correlations, on the other hand, showed
23 consistently low agreement with the simulated transcript groups. Finally, we considered the

1 number of transcripts that remained unclustered (Figure 4C) before and after re-clustering
2 unassigned transcripts (i.e. cluster expansion step, see Methods). Pearson correlation
3 provided successful cluster assignment for practically all transcripts in the simulated dataset,
4 especially when incorporating percentiles (Pearson: ~10% unclustered before expansion, ~1%
5 after; percentile: ~4% unclustered before expansion, 0% after), whilst the rest of metrics
6 performed significantly worse, leaving 20-30% of transcripts unassigned even after cluster
7 expansion, with proportionality (~30% unclustered before expansion, ~25% after) being the
8 less optimal. Altogether, though ρ demonstrated good performance in many aspects of
9 clustering, including intra-cluster correlation and agreement with the simulated clustering, it
10 was outperformed by percentile correlation when globally considering all evaluated
11 parameters (Figure 4D). In addition to the fact that ρ failed to control for unassigned
12 transcripts, computing means and standard deviations of Jaccard indices across simulated-
13 calculated pairs showed percentile and Pearson correlations as the most consistently
14 accurate methods. All in all, our synthetic data evaluations showed that the percentile
15 correlation approach performed well -and more consistently than ρ proportionality- in all the
16 evaluated features, and visibly captured co-expression better than both traditional and zero
17 inflation-aware correlation metrics.

18 *Co-Differential Isoform Usage analysis of single-cell isoform expression*

19 Isoform clusters represent groups of alternative transcripts that are co-expressed at the cell
20 type level. However, our clustering results did not provide information on iso-transcriptome
21 properties associated with splicing regulation. To facilitate interpretation of the isoform

1 clustering results, we first defined genes with Differential Isoform Usage (DIU) as those whose
2 isoforms were assigned to different clusters (Figure 5A) and found that 80% of genes with
3 clustered isoforms (2,577 out of 3,172) were DIU and involved a total of 7,575 isoforms. In this
4 context, DIU genes will necessarily have two or more isoforms with significant changes in
5 expression across cell types and simultaneously undergo cell type-dependent post-
6 transcriptional regulation, leading to changes in isoform expression in each cell type.

7 In order to study isoform co-expression patterns, we defined co-Differential Isoform Usage
8 (coDIU) genes as those showing coordinated cell-type-specific isoform usage. Specifically, we
9 considered two or more genes to be coDIU if their isoforms had been assigned to the same
10 clusters (Figure 5B, Methods). This resulted in the definition of an isoform co-splicing network
11 where nodes are clusters of correlated isoforms and edges represent coDIU genes, i.e. the
12 number of genes for which two or more isoforms are co-expressed across cell types (Figure
13 5C, Methods). To ensure the selection of significantly coDIU genes, we fitted a generalized
14 linear model for every pair of coDIU genes, and selected pairs with isoform-level co-
15 expression and no significant cell type expression variation when only accounting for gene-
16 level expression (see Methods). CoDIU genes therefore present cell type-dependent co-
17 expression of at least two isoforms, represented by cluster assignment matches, but are not
18 co-expressed when only gene expression is considered. Using this strategy, we detected 2,049
19 genes with at least one significant coDIU partner ($cluster*cell-type$ FDR < 0.05, $gene*cell-type$
20 FDR > 0.05), involving 6,370 co-expressed isoforms. The number of coDIU genes sharing
21 isoforms across each cluster pair was variable, although it rose up to ~120 for highly
22 connected expression profiles (Figure 5C).

1 We then interrogated the coDIU network to find patterns underlying the splicing coordination
2 signal detected by the *acorde* pipeline. First, in order to measure whether coDIU generated
3 strong or subtle variations in isoform selection across cell types, we investigated the
4 association of coDIU to single or multiple cell type isoform switching events. Note that single
5 isoform switching events involve clusters with patterns that are similar across all cell types
6 except one, leading to high between-cluster correlations. Interestingly, we found that the
7 number of coDIU genes linking isoform co-expression clusters was dependent on cluster
8 sizes, but showed no direct relationship with the similarities between expression profiles
9 (Figure 5D). The detection of coDIU genes involving isoforms with highly different expression
10 patterns suggested that coordinated isoform usage mechanisms are able produce strong cell
11 type-level shifts in isoform selection, and thus modulate the expression of highly cell type-
12 specific splice variants.

13 We next evaluated the cell type-level relationships present in the isoform co-expression
14 network, namely the occurrence of coDIU across all possible pairs of cell types in our data.
15 Although co-splicing could potentially occur between any combination of cell types, our
16 results showed that a high proportion of coDIU interactions were detected when the isoforms
17 involved had high expression in one of the two neural cell types, i.e. GABA-ergic and
18 glutamatergic neurons (Figure 5E). This can be partially explained by the fact that some of the
19 clusters with neuron expression are among the largest generated by the *acorde* pipeline
20 (Figure 5C). However, another plausible explanation is that the central role of neurons in the
21 tissue under study (i.e. primary visual cortex) might situate co-splicing at the core of neural
22 function regulation, as well as the modulation of its interaction with glial cell types.

1 *Functional analysis of the coDIU network*

2 We next set out to investigate the functional implications of our isoform co-expression
3 network. Since, with a few exceptions^{59,60}, splicing analysis tools rarely integrate functional
4 information and function annotation databases do not usually include data at isoform
5 resolution, we annotated the long read-defined transcripts using IsoAnnotLite
6 (<https://isoannot.tappas.org/isoannot-lite/>), included both transcript and protein-level motifs,
7 sites and domains, as well as non-positional, gene-level features such as Gene Ontology (GO)
8 terms (for a detailed description of the annotation process and a comprehensive list of
9 functional categories and source databases, see [Supplementary Note 1](#)).

10 First, we analyzed which biological processes and gene functions were potentially controlled
11 by DIU (AS-regulated) and coDIU (co-regulated) mechanisms operating across cell types, i.e.
12 which functions were overrepresented at DIU and coDIU genes. In order to discriminate the
13 functional properties of AS-regulated genes from those showing no cell type specificity in
14 isoform expression, we performed a functional enrichment test of DIU genes vs genes with DE
15 isoforms, which were used as the background ([Figure 6A](#), see Methods). Interestingly, DIU
16 genes showed significant enrichment (FDR < 0.05) of GO terms associated with transcription
17 (*gene expression, DNA-templated transcription*) and RNA metabolism (*RNA processing, RNA*
18 *metabolic process*), as well as protein features required for these processes, such as Nuclear
19 Localization Signals (NLS), Post-Translational Modifications (PTMs), particularly acetylation
20 sites (Acetyl-Lys), and protein-complex formation ([Figure 6A](#)). In order to investigate the
21 cellular processes where coDIU could potentially have a relevant regulatory role, we

1 compared the proportion of coDIU and DIU genes annotated for each functional feature in
2 the transcriptome using a partially-overlapping samples z test⁶¹ (Figure 6B, see Methods). Not
3 surprisingly, coDIU genes were significantly enriched (FDR<0.05) in some of the same
4 functionalities obtained in the previous DIU analysis, such as RNA biosynthesis and processing
5 (*transcription, gene expression, RNA metabolic process*). However -and in contrast to DIU genes,
6 which were markedly involved in biosynthesis-, genes regulated by coDIU were enriched in
7 mitochondria components (i.e. *mitochondrial matrix*), suggesting that coordinated isoform
8 usage may affect oxidatoin and energy metabolism. Remarkably, coDIU genes also showed
9 additional enrichment for splicing-related terms such as *RNA splicing, mRNA splicing via*
10 *spliceosome* and for *3' UTR motif K-box* (Figure 6B). This result links genes involved in splicing
11 and elements affecting RNA stability with the coordination of AS and suggests that the co-
12 expression of alternative isoforms is a post-transcriptionally regulated process, in line with
13 previous evidence of AS self-regulation⁶²⁻⁶⁵.

14 To obtain further insight into the functional elements controlled by AS co-regulation across
15 neural cell types, we performed a Functional Diversity Analysis (FDA, Figure 6C). FDA is part of
16 the tappAS framework⁶⁰, and identifies functionally *varying* genes, i.e. genes expressing
17 transcript variants with differences in the inclusion of functional features (see Methods). FDA
18 can be evaluated from a presence/absence standpoint (i.e. AS completely removes a feature),
19 or by detecting variation in the transcript positions defining the feature (Figure 6C). We
20 therefore compared the diversity in transcript-level functional features between DIU, coDIU
21 genes and genes with more than one DE isoform, for all functional categories provided by
22 IsoAnnotLite (Figure 6D). Interestingly, we detected a larger percentage of varying genes as
23 the level of considered regulatory complexity increased, with coDIU resulting in the largest

1 amount of feature inclusion diversity in virtually all protein and transcript feature categories.
2 To measure this effect, we compared varying percentages between all combinations of the
3 three gene sets (paired samples t-test, see Methods) and confirmed the observed trend,
4 regardless of whether the variability criteria employed was position or presence. In particular,
5 even though all comparisons were significant, coDIU resulted in the most significant increase
6 in feature variation (coDIU vs DE isoform genes p -value: presence = 7.61×10^{-7} , position = 3.45×10^{-5} ;
7 coDIU vs DIU genes p -value: presence = 4.92×10^{-4} , position = 9.62×10^{-3}). We verified that
8 this increase in functional diversity was not associated to expression level or isoform length
9 biases in coDIU genes (Supplementary Figures 14A-B). This result suggests that alternative
10 isoforms that engage in co-expression relationships tend to alter their functional properties
11 significantly more often than other transcripts, namely by controlling the inclusion of motifs,
12 sites and /or domains, thereby coupling alternative splicing and isoform co-expression with
13 functional potential.

14 *Functional analysis of neuron-oligodendrocyte isoform co-expression*

15 To further understand the relationship between cell-type identity, isoform co-expression and
16 the functional properties of coDIU genes, we searched the coDIU network for cluster groups
17 representing biologically-related isoform switches between defined neural types. Namely, we
18 focused on a set of genes 160 coDIU genes (Figure 5C) representing oligodendrocyte-specific,
19 neuron-specific, and shared isoform expression patterns (clusters 8, 13 and 19, respectively,
20 Figure 6E) and analyzed isoform-associated functional variability using FDA. For this set of
21 alternative isoforms, polyA site and 3'UTR length showed the highest variation rates among

1 annotated transcript-level functional categories (varying in $\geq 70\%$ genes, [Supplementary Figure](#)
2 [15A](#)). Moreover, we noticed that these changes followed a clear cell type-specific pattern, with
3 the majority of coDIU genes shared among the three clusters expressing their longest 3'UTR
4 isoforms in neurons ([Supplementary Figure 15B](#)) and neural-specific isoforms generally
5 expressing longer 3'UTRs than their oligodendrocyte-specific counterparts ([Supplementary](#)
6 [Figure 15C](#)).

7 Next, we inspected 3'UTR-related functional categories, i.e. RBP binding, miRNA binding, and
8 3'UTR motifs, using ID-level FDA (see Methods) to identify specific functional features
9 associated to isoform usage differences between neurons and oligodendrocytes. Regarding
10 the presence of RBP binding sites, we found that Mbnl and CELF4 binding sites showed high
11 variation levels, i.e. 50% and 75% of genes, respectively ([Supplementary Figure 15D](#)); however,
12 the amount of genes including Mbnl and CELF4 binding motifs was low ($\sim 3\%$ genes), possibly
13 due to the relatively poor annotation density of RBP sites transferred by IsoAnnotLite
14 ([Supplementary Note 1](#)). We also detected high variation frequencies for several miRNA sites
15 ([Supplementary Figure 15D](#)), although no specific miRNA motif was shared by more than $\sim 5\%$
16 of genes (up to 9 out of 160 genes). Nevertheless, regarding 3'UTR motifs, we found that the *K-*
17 *box* motif presented inclusion changes in $\sim 60\%$ of annotated coDIU genes ([Figure 6A](#)).

18 Interestingly, *K-box* motifs have been described to have a role in negative post-transcriptional
19 regulation by promoting miRNA binding and transcript degradation^{66,67}. In line with this, the
20 coDIU network included several genes in which 3'UTR elongation led to neuron-specific co-
21 inclusion of *K-box* motifs, included *Atxn10* (ataxin 10, [Supplementary Figure 16A](#)) and *Btrc*
22 (beta-transducin repeat containing gene, [Supplementary Figure 16B](#)), both of which have been
23 proposed to be involved in neuron survival and differentiation^{68,69}. These results suggest that

1 a 3' UTR binding-mediated mechanism favoring isoform co-expression may operate to fine-
2 tune the post-transcriptional regulation of neuron survival genes.

3 Importantly, the majority of neuron-oligodendrocyte coDIU genes also presented coding
4 region variation (i.e. CDS, [Supplementary Figure 15A](#)), revealing that the coordination of
5 isoform usage can modify both transcript and protein functional properties. In particular,
6 protein domains (PFAM) and post-translational modifications (PTMs) presented high variation
7 rates (varying in ~40% of genes, [Supplementary Figure 15A](#)) and thus constituted the
8 categories with the most cell type-level dependent functional variation. While ID-FDA reported
9 no specific PFAM domains shared among the analyzed gene set (~1%, maximum of 2 out of
10 160), up to ~14% of them presented inclusion variation in similar PTMs (23 out of 160) with
11 medium to high variation rates for phosphorylation, acetylation, ubiquitination and ligand
12 binding sites ([Supplementary Figure 15D](#)). However, synergies between PTM and domain
13 inclusion changes could still result in differential functional activities at the protein level. As an
14 example, we found two tubulin isotypes, i.e. γ -tubulin 2 (*Tubg2*) and β -4-tubulin B-chain
15 (*Tubb4b*), that co-expressed isoforms in clusters 8 and 19, that is, protein variants with
16 neuron-specific and neuron-oligodendrocyte expression, respectively. Interestingly, both of
17 these genes also presented inclusion changes in an N-terminal GTP-ase protein domain and
18 several PTMs ([Figure 6F](#) and [Figure 6G](#)), although with different functional outcome. *Tubg2*
19 presented neuron-specific expression of the full-domain, PTM-including isoform, resulting in
20 the inclusion of a phosphoserine residue and a GTP binding motif ([Figure 6F](#)) and suggesting
21 the GTP-ase role to be enhanced in mouse visual cortex neurons in comparison to
22 oligodendrocytes. On the other hand, the domain-including *Tubb4b* isoform was expressed in
23 both cell types ([Figure 6G](#)), suggesting that this β -tubulin isotype has a broader GTP-ase

1 function than *Tubg2*. The domain/motif inclusion pattern arising from isoform co-expression
2 in these two tubulin isotypes may reflect the different cellular roles of γ and α/β tubulins⁷⁰.
3 Specifically, α/β tubulins are the building-blocks of the microtubule cytoskeleton, while γ -
4 tubulin has a major role in microtubule nucleation and on the regulation of microtubule
5 dynamics. The γ -tubulins only present two isotypes, among which *Tubg2* is exclusive to
6 neurons and has a well-defined role in neuron survival and growth⁷¹, whereas *Tubb4b* is part
7 of a broader catalogue of β -tubulin isotypes, all of which are jointly regulated and share a
8 structural function⁷². Intriguingly, however, mutations in a closely-related isotype (*Tubb4a*)
9 have already been shown to have a differential phenotypic effect in neurons and
10 oligodendrocytes⁷³. In line with this, we observed modifications in the number and density of
11 N-terminal phosphorylation sites in *Tubb4b* as a result of cell type-level isoform co-expression,
12 with an increase in protein diversity in neurons (3 vs 2 alternative isoforms, **Figure 16G**). Even
13 more interestingly, PTM changes on different α/β -tubulin isotypes have long been known to
14 modify tubulin stability and its interactions with other proteins⁷⁴, creating a “tubulin code”⁷⁵.
15 As a result, tubulin PTMs have the ability to modulate microtubule stability and function,
16 intervening in processes such as neural differentiation, survival and polarity^{76,77}. This cell type-
17 specific inclusion of different PTMs via isoform co-expression could therefore operate as a
18 fine-tuning mechanism of the tubulin code, modifying the number and position of available
19 modification sites in β -4-tubulin.

20 **Discussion**

1 Alternative Splicing (AS) is known to be a tightly-regulated process in which splicing factors
2 interact to create cell type-specific isoform expression patterns⁷⁸. The transcriptome-level
3 consequences of AS regulation have been studied in different ways, including, but not limited
4 to, the detection of within-isoform coordination of alternative sites^{40,41}, the generation of
5 gene-isoform networks to uncover novel regulatory relationships⁷⁹⁻⁸² and the application of
6 single-cell data to unravel cell type-specific expression patterns for same-gene isoforms^{27,83}.
7 However, the extent to which AS regulation creates co-expression patterns among alternative
8 isoforms from different genes has not yet been fully addressed. Specifically, previous studies
9 tackling this type of isoform co-expression have either focused on specific event types, such
10 as alternative 3' exons⁸⁴, or solely on the identification of functionally-relevant alternative
11 isoforms in different biological contexts^{85,86}.

12 In this study we present *acorde*, an end-to-end pipeline to generate isoform co-expression
13 networks and detect genes with co-Differential Isoform Usage (coDIU), and apply it to the
14 study of isoform co-expression among seven neural broad cell types⁴³. To this end, we
15 successfully leveraged single-cell data by implementing percentile correlations, a metric
16 designed to overcome single-cell noise and sparsity and provide high-confidence estimates of
17 isoform-to-isoform correlation. Here, we show that percentile-summarized Pearson
18 correlations outperform both classic and single-cell specific correlation strategies⁵⁶, including
19 proportionality methods that were recently proposed as one of the best alternatives to
20 measure co-expression in single-cell data⁴². In addition, using a long read-defined, functionally
21 annotated transcriptome enabled us to obtain a biological readout from the isoform network.
22 coDIU genes were found to be enriched in the same biological functions, a number of which
23 were unique in comparison to enrichment results for genes solely reported as DIU. Inter-gene

1 isoform co-expression thus appears to impact a subset of DIU genes sharing a specific set
2 functions, which suggests that coDIU may contribute with an additional layer of complexity to
3 some cellular processes, operating as a fine-tuning mechanism. Our analyses also revealed
4 that isoforms from coDIU genes encompass higher functional diversity than those belonging
5 to DIU genes, an effect that was not associated to expression or length differences.
6 Importantly, these changes impact both transcript and protein isoform functional features,
7 which pinpoints their ability to globally increase the functional repertoire of coDIU genes.
8 Mechanisms generating isoform co-expression can therefore be thought of as a potential
9 source of functional synergies between alternatively spliced genes, giving rise to simultaneous
10 changes in functional properties among co-expressed isoforms.

11 To demonstrate the power of the *acorde* pipeline, we include examples where these kinds of
12 coordinated changes were detected in the mouse neural dataset. First, we report a neural-
13 specific pattern of 3'UTR co-elongation that is consistent with the available literature^{87,88} and
14 results in a simultaneous increase in the number of K-box motifs, miRNA and RBP binding
15 sites in these UTRs. In addition, we describe how cell type-specific co-expression of γ and β -
16 tubulin isoforms creates divergent functional properties among them, which is consistent with
17 their different biological functions^{71,72} and points to a splicing-mediated coordination of cell-
18 type specific structural components. On a broader note, an interesting finding stemming from
19 our functional analyses concerns the type of biological properties that make the functional
20 signal of the coDIU network. Namely, while we do find enrichment of biological processes (i.e.
21 *RNA metabolism* or *mRNA splicing via spliceosome*) and differential inclusion of many features
22 from a wide variety of functional categories, we failed to recover specific functional elements
23 consistently present in a large number of co-expressed isoforms, such as binding sites for

1 specific RBPs or microRNAs. This suggests that AS and post-transcriptional regulation
2 mechanisms operate by controlling the presence of a broad array of functional features,
3 which jointly contribute to modulate key cell-level processes encoding cell-type identity. While
4 these insights need to be subject to further experimental validation, they serve to illustrate
5 the hypothesis-generating power of our pipeline.

6 All in all, we have hereby showed that *acorde* can effectively leverage single-cell RNA-seq data
7 to build isoform co-expression networks, revealing a new dimension of post-transcriptional
8 regulation of gene expression while also disclosing the cellular processes and functional
9 elements impacted by these mechanisms.

1 **Methods**

2 **Single-cell data pre-processing and quality control**

3 Mouse neural single-cell RNA-Seq data from mouse primary visual cortex was obtained from
4 Tasic *et al.*⁴³ and consists in paired-end Illumina reads generated with the Smart-seq2
5 protocol⁸⁹, which enables isoform-level quantification. Reads were downloaded from
6 Sequence Read Archive accession SRP061902 and mapped to the mouse genome
7 (GRCm38.p6) using STAR⁹⁰. We performed isoform expression quantification of the long read-
8 defined isoforms (see [Supplementary Note 1](#) for details on long read transcriptome definition)
9 using RSEM⁹¹, and used the labels provided by Tasic *et al.* to assign the 1679 cells to 7 broad
10 cell types: 5 glial (microglia, endothelial cells, oligodendrocytes, oligodendrocyte precursor
11 cells (OPCs) and astrocytes) and 2 neural (GABA-ergic and glutamatergic neurons).
12 Isoform length effect on expression was evaluated using the *NOISeq* R package⁹², where mean
13 expression showed to be highly correlated with transcript length (adjusted $R^2 = 0.81$; p-value =
14 $2.2e-16$). Using isoform i effective length (l_i) and cell-level j estimated counts (c_{ij}), both output
15 by RSEM and, after testing several alternatives, we devised a custom formula to minimize the
16 impact of length on isoform expression for each isoform i :

17
$$y_{ij} = \frac{c_{ij}}{\left(10^{-6} \sum_{i=1}^I c_{ij}\right) \sqrt{10^{-3} l_i}} \quad \text{[Equation 1]}$$

18 The transformed expression value for isoform i in cell j (y_{ij}) was again tested for length bias
19 and a low correlation was found (adjusted $R^2 = 0.25$; p-value = $6.84e-8$). Next, we inspected the
20 library size distribution and filtered both high and low-count outliers due to potential

1 premature cell death or library preparation duplets, with a total of 1591 cells passing quality
2 control. Feature-level quality control was performed in a cell type-aware manner, keeping
3 isoforms that showed non-zero expression in at least 20% of one cell type. Out of the 36,986
4 isoforms and 12,692 genes in the PacBio-defined transcriptome, we retained 18,867 isoforms
5 and 9,596 genes for downstream analysis. Finally, we retained isoforms from multiple-isoform
6 genes, hence removing cases where no AS-directed, isoform-level co-expression relationships
7 can be established. This sets a general requirement for the entire study, which is that all
8 isoforms retained must have, at all times, at least one same-gene counterpart to establish
9 regulatory relationships that can be based on differential splicing of that gene, given that no
10 AS regulation can be detected if a gene's total expression is represented by a single isoform.
11 As a result, 13,832 isoforms from 4,591 genes were preserved.

12 **Differential Expression (DE) across multiple groups**

13 *Single-cell DE analysis*

14 Differential Expression (DE) analysis among the 7 cell types was performed by combining
15 ZinBWaVE weights⁴⁷ and bulk-designed DE methods edgeR⁴⁹ and DESeq2⁴⁸ (i.e. using the
16 corresponding R packages), which enable multiple group testing and were among the best-
17 performing methods when combined with the ZinBWaVE method. Briefly, ZinBWaVE
18 calculates cell-level weights for each isoform, effectively downweighting zeros during
19 modelling for differential expression in single cell data (see van den Berge *et al.*⁴⁷ for details),
20 and hence unlocking bulk RNA-Seq computational methods for single-cell data. Of note,
21 generalized linear models (GLM) within edgeR and DESeq2 were built and run following the

1 pipeline used by van den Berge et al.⁴⁷ to make them suitable for single-cell RNA-seq data, and
2 are implemented in a wrapper function within our R package as described in Equation 2,
3 where y_{ij} is expression of isoform i in cell j , T_{kj} is a dummy variable which takes value 1 when
4 cell j is assigned to cell type k ($k=1, \dots, K$) and 0 otherwise, β_{ki} are the regression coefficients for
5 isoform i , ε_{ij} represents the error term, and $h()$ is the link function of the GLM (natural
6 logarithm in this case).

7
$$h(y_{ij}) = \beta_{0i} + \sum_{k=2}^K \beta_{ki} T_{kj} + \varepsilon_{ij} \quad \text{[Equation 2]}$$

8 Differential Expression was defined using a significance threshold of $FDR < 0.05$ when testing
9 the significance of the model for each isoform i , that is, $H_0: \beta_{2i} = \dots = \beta_{Ki}$. Isoforms considered
10 DE were preserved for downstream analysis if detected by at least one of the two methods
11 edgeR or DESeq2, since this indicates a change in expression for any of the cell types
12 considered rather than a flat expression profile. In addition, isoforms were removed if they
13 belonged to genes with a single DE isoform, in agreement with the idea that no differential
14 splicing regulation can be detected for single isoform genes, and hence that they will not form
15 co-splicing relationships with isoforms from other genes in the network.

16 Prior to DE testing, and to balance the number of cells per cell type, the most abundant cell
17 types (GABA-ergic neurons, $n = 761$; glutamatergic neurons, $n = 764$ cells) were downsampled
18 by randomly selecting 45 cells, keeping $N=241$ cells for downstream analysis, simultaneously
19 reducing the computational cost of multi-group DE testing. Finally, since downsampling can
20 lead to having isoforms with all-zero counts again, we re-filtered isoforms that retained non-

1 zero expression in at least 20% of one cell type, ultimately keeping 13,451 isoforms from
2 4,583 genes for the DE analysis.

3 *Validation of downsampling and single-cell DE analysis*

4 To check that the random selection of cells in the downsampling process was not affecting the
5 DE results, we performed 50 independent runs of random neural cell sampling (without
6 replacement) followed by zero-expression filtering (expression above zero for at least 25% of
7 cells in at least one cell type, where threshold stringency was increased with respect to the
8 standard pipeline to best control zero abundance among iterations and avoid problems
9 during GLM modelling) and DE testing with edgeR and DESeq2. First, to measure the
10 consistency of each method independently, we calculated the mean and standard deviation of
11 the number of DE isoforms across all same-method sampling runs ($R = 50$). To check the level
12 of within-method agreement, we next considered isoform IDs labeled as DE in each
13 independent method run, and calculated the Jaccard Index (J_{rs}) between DE results of that
14 same method for all possible pairs of random sampling runs r and s ($r, s = 1, \dots, 50$, $r < s$, a total of
15 1225 comparisons). To summarize this information, we relied on the mean and standard
16 deviation of these two sets of J_{rs} values. Finally, we measured the level of agreement between
17 edgeR and DESeq2 regarding our DE criteria, that is, considering isoforms detected by at least
18 one of the methods to be significantly DE ($FDR < 0.05$). To achieve this, we calculated the union
19 of DE isoforms between one-to-one pairs of edgeR and DESeq2 runs ($R = 50$), and computed
20 the Jaccard Index between all possible pairwise combination of global DE results, i.e. isoforms
21 detected by at least one method (again, 1,225 comparisons).

1 **Percentile correlation**

2 In order to assess the similarity of isoform expression profiles across cells, a correlation
3 measurement can be used by taking cells as observations. We propose here instead to first
4 summarize the expression within a given cell type with percentiles and then compute the
5 correlation using all cell types and their percentiles as observations, a method that we refer to
6 as “percentile correlation” (see main text [Figure 2A](#)).

7 Percentile correlations rely on the assumption that cell-to-cell differences can be mostly
8 attributed to transcriptional stochasticity or technical noise, and that these within-cell type
9 differences have a smaller effect than between-cell type expression differences. However,
10 expression estimates for transcripts within the same cell are biased in different degrees,
11 mostly depending on their expression levels, with lower expression being generally
12 accompanied by higher noise levels⁵². This modifies the extent to which isoforms are affected
13 by noise in each cell and causes strong cell-level effects that prevent the detection of co-
14 expression relationships using solely cell-level measurements. Instead, we set out to target
15 changes in expression across cell groups. We therefore considered isoform expression levels
16 in the different cell types as a range of possible values, defined by the cell-level
17 measurements in the data. In this context, the expression value of an isoform in a cell is used
18 a proxy to infer the underlying distribution of expression values in the cell type, where the
19 shape and width of this distribution will depend on both biological and technical factors.

20 To translate this into a metric, we first took the expression values of an isoform in each of the
21 cell types and computed a number of percentiles (p). We selected $p = 10$ to achieve a good
22 balance between accuracy and computational burden in downstream analysis. As the

1 minimum expression value (percentile 0) was also included, we actually had 11 new values
2 representing the expression range within a given cell type. As a result, each isoform will
3 possess a new, recalculated expression vector where the percentile values computed in each
4 cell type will replace cell-level expression estimates. This process was repeated for each
5 isoform. Next, we computed pairwise Pearson correlations between every pair of isoforms,
6 obtaining a percentile correlation matrix \mathbf{R} . In this context, high correlations will appear if a
7 pair of isoforms shows a similarly broad expression distribution in most cell types, as well as a
8 similar amount of relative expression change between cell types.

9 **Semi-automated isoform clustering**

10 In order to obtain modules of tightly co-expressed isoforms, we combined the hierarchical
11 clustering algorithm with several rounds of cluster profile refinement (see main text, [Figure](#)
12 [3A](#)), in order to automate the most intensive steps of clustering while also granting control
13 over the level of aggregation and within-cluster similarity. Clustering and refinement steps can
14 be combined and re-arranged to best capture co-expression patterns within the data, and
15 their parameters can be defined by potential future users to provide maximum flexibility.
16 Functions for clustering and refinement are implemented in the *acorde* R package.

17 *Dynamic hierarchical clustering*

18 The previously obtained correlation matrix (\mathbf{R}), where each element r_{ij} represents the
19 Pearson's correlation coefficient between the percentiles of isoforms (i,j), was transformed
20 into a distance metric to be used in the hierarchical clustering. As we aimed to cluster

1 positively correlated isoforms given our biological hypothesis, we discarded negative
2 correlation values by replacing them with zero values, and therefore defined the distance
3 between any pair of isoforms i and j as in Equation 3.

$$4 \quad d_{ij} = \begin{cases} 1 - r_{ij} & \text{if } r_{ij} > 0 \\ 1 & \text{if } r_{ij} \leq 0 \end{cases} \quad \text{[Equation 3]}$$

5 We next performed hierarchical cluster analysis using the *hclust()* function in the R *base*
6 package⁹³ with the average linkage criterion, and obtained a dendrogram. To obtain clusters,
7 we used the *cutreeHybrid()* function in the *dynamicTreeCut* R package⁵⁵ in order to find different
8 thresholds for different branches of the dendrogram tree, instead of using a fixed threshold
9 for the entire dendrogram. We provided the following non-default parameters to the
10 *cutreeHybrid()* function: *deepSplit* = 4, *pamStage* = FALSE, *minClusterSize* = 20. Briefly, *deepSplit*
11 ranges between 0 and 4, and provides smaller clusters, more accurate clusters when set to
12 high values. *pamStage* determines whether a second stage of clustering using an algorithm
13 similar to the Partition Around Medoids (PAM) method will be performed after searching the
14 dendrogram for clusters (see Langfelder *et al.*⁵⁵). As a result of this PAM-like step, no items are
15 left unassigned to clusters, while setting *pamStage* to FALSE allows unclustered items. Finally,
16 *minClusterSize* determines the minimum size of the produced clusters, and thus passing a
17 higher value to this argument prevents the generation of too many clusters with a very small
18 number of items.

19 This initial set of clusters is to be used as “hooks” to gather as much expression profile
20 diversity from the data as possible. Importantly, even though our parametrization allows
21 isoforms to remain unassigned to clusters (see above), some isoforms may still show low

1 similarity to their cluster's profile. To be able to obtain profiles as consistent as possible for
2 downstream refinement, we implemented a cluster quality control step in which we removed
3 isoforms based on a minimum correlation threshold with the rest of the members. In
4 particular, isoforms were moved to the unclustered group if they showed a correlation lower
5 than 0.85 with 3 or more isoforms from their cluster. In this manner, only tightly-correlated
6 groups of isoforms will remain clustered.

7 *Expanding clusters with unassigned isoforms*

8 To re-assign unclustered isoforms to clusters with which they show high correlation, we
9 performed two rounds of correlation-based, cluster expansion. In this process, we first
10 summarized each cluster's profile into a synthetic representative transcript that we named
11 "metatranscript". Metatranscripts were calculated as the mean of the percentile-based
12 expression of all isoforms in the cluster. As a result, we obtained $11 \cdot K$ (K being the number of
13 cell types) mean-summarized percentile expression values, which can be understood as an
14 approximation to the expression range shown by the isoforms from that cluster in each of the
15 cell types.

16 We next computed correlations between metatranscripts and unclustered isoforms, and
17 performed two rounds of assignment. First, we assigned unclustered isoforms showing
18 correlation values above 0.9 with at least one cluster, where the maximally correlated cluster
19 was selected as the best match if there were ties. Next, metatranscripts were re-calculated for
20 the newly expanded clusters and a similar assignment round was performed, this time
21 lowering the correlation threshold to 0.8. In doing this, unclustered isoform groups are

1 assigned in order, and highly correlated elements can therefore contribute to strengthen
2 within-cluster similarities before assigning more lowly correlated elements.

3 *Merging clusters by profile similarity*

4 Prioritizing the reduction of within-cluster variability led to obtaining a large number of small,
5 redundant clusters. To mitigate this effect while also preserving high correlations between
6 cluster members, we merged clusters by profile similarity using the correlations between their
7 metatranscripts. To achieve this, we performed dynamic hierarchical clustering using the
8 *cutreeHybrid()* function on metatranscripts, with the following non-default parameters:
9 *deepSplit = 4*, *pamStage = FALSE* (see section on dynamic hierarchical clustering above for more
10 information), *minClusterSize = 1*. Of note, *minClusterSize* is set to 1 to avoid forcing merge of
11 clusters that do not show enough similarity. As a result, we joined highly similar clusters,
12 decreasing the number of clusters and hence minimizing redundancy.

13 While this strategy effectively eliminated between-cluster redundancy, we observed some
14 exceptions where clearly similar profiles were not merged by clustering metatranscripts
15 ([Supplementary Figure 3A](#)). In addition, a few cases arose where automated merging resulted
16 in joining clusters with highly uncorrelated profiles ([Supplementary Figure 3B](#)). Given that it is
17 very hard to provide a completely automated solution to the clustering problem, we solved
18 these final discrepancies using a series of manually supervised refinement steps.

19 For that, we first inspected cluster profiles to flag incorrect merge decisions generating noisy
20 clusters. Isoforms in these clusters were treated as unclustered, and used to expand the
21 remaining, well-defined clusters. To do this, we computed the correlation between the

1 isoforms to be assigned and cluster metatranscripts, similarly to the cluster expansion
2 process described above. In this case, we joined isoforms to clusters if they presented
3 Pearson correlation > 0.8 with its metatranscript. Isoforms showing high correlation with
4 several clusters were assigned to the cluster with maximum correlation. Once isoforms from
5 noisy clusters had been re-assigned, clusters were manually inspected again to detect profile
6 redundancies missed by automated metatranscript clustering. Clusters that presented a
7 similar expression pattern across cell types were then merged.

8 **Co-expression pattern simulation**

9 To validate percentile correlations and our clustering strategy, we evaluated their
10 performance on synthetic data, where co-expression relationships between simulated
11 features need to be pre-defined as part of the data simulation process. However, there is, to
12 the best of our knowledge, no currently available strategy to simulate single-cell data
13 including modules of co-expressed features. We therefore designed our own simulation
14 strategy by combining the *SymSim* R package⁵⁸ to adequately model single-cell RNA-Seq data,
15 and a dedicated strategy to generate co-expression between *SymSim* simulated features.

16 First, we set the following parameters to the *SimulateTrueCounts()* function in *SymSim* in order
17 to obtain a count matrix consisting in 1000 cells from 8 cell types and 8000 features, with
18 sufficient feature-level variation between the different cell groups:

```
19 SimulateTrueCounts(ncells_total = 1000, min_popsize = 100, i_minpop  
20 = 1, ngenes = 8000, nevtf = 10, n_de_evtf = 9, evf_type = "discrete",
```

```
1     phyla = pbtree(n = 7, type = "discrete", vary = "s", Sigma = 0.25,  
2     gene_effect_prob = 0.5, bimod = 0.4, prop_hge = 0.03, mean_hge = 5)
```

3 Next, we modeled technical effects on these true counts in order to obtain real, observed
4 counts using the *True2ObservedCounts()* function in *SymSim*, with the following parameters:

```
5     True2ObservedCounts(true_counts$counts,  
6     meta_cell = true_counts$cell_meta, protocol = "nonUMI",  
7     alpha_mean = 0.1, alpha_sd = 0.005, lenslope = 0,  
8     gene_len = rep(1000, nrow(true_counts$counts)), depth_mean = 4e6,  
9     depth_sd = 1e4)
```

10 To create co-expression patterns, we then re-ranked expression values on a cell type-specific
11 manner to define synthetic features, based on the expression profile of 15 pre-defined co-
12 expression modules.

13 First, we drafted 15 different co-expression profiles reflecting three levels of expression
14 complexity, that is, showing high expression or expression "peaks" in one, two, or three cell
15 groups, respectively. To generate a count matrix reflecting these expression patterns, we
16 shuffled simulated counts to create new, synthetic features. To achieve this, we first re-ranked
17 features in each cell group by mean expression across cells in the group, breaking feature
18 connectivity between the simulated cell types. Then, the top 1,400 features from each cell
19 type were selected, together with the bottom 1,400 features. In this manner, we obtained
20 high-expression and low-expression count vectors for each group, which we then combined
21 to create synthetic features following the pre-designed cluster's co-expression pattern. For
22 each cluster, 200 count vectors from top-expression features were assigned to peaking
23 groups, and 200 count vectors from bottom-expression features to cell groups showing low

1 expression. Of note, 1,400 features were selected for simulation in order to grant at least 7
2 different 200-feature groups could be generated for each cell type, where the cell group with
3 the highest peaking frequency across clusters showed high expression in 6 clusters only.

4 All in all, we obtained a simulated count matrix containing 1,000 cells from 8 cell types and
5 3,000 synthetic features, all of which belong to one of the 15 simulated co-expression
6 modules. Therefore, by breaking feature-level connectivity between cell types, we benefited
7 from feature-specific properties at the cell type level, while re-creating cell type expression
8 coordination patterns that the SymSim strategy was not able to generate. Finally, to ensure
9 the quality of the simulated clusters, we filtered synthetic features if their Pearson correlation
10 with the cluster's median profile was below 0.75 ([Supplementary Figure 4B](#)).

11 **Benchmarking of isoform correlation metrics for scRNA-seq data**

12 Traditional correlation metrics have been shown to perform poorly when applied to scRNA-
13 seq data, mainly given the increased noise and stochasticity levels in this data type. Recently,
14 extensive benchmarks including single cell-tailored metrics have shed light on how to best
15 select correlation metrics for single-cell data (see review by Skinnider *et al.*⁴²). We therefore
16 compared the performance of percentile correlations to a representative set of correlation
17 metrics used in single-cell co-expression studies, namely classical Pearson and Spearman
18 correlations, single-cell designed zero-inflated Kendall⁵⁶ correlation, and proportionality
19 metric rho (ρ)⁵⁷, in agreement with previous reports showing that proportionality metrics were
20 among the best performing co-expression methods in single-cell data. To measure

1 performance, we computed these five co-expression metrics for all the synthetic features in
2 the previously-simulated dataset, generating five different distance matrices for clustering,
3 and evaluated which metric best recapitulated the simulated co-expression modules when
4 used in our clustering pipeline. Pearson and Spearman correlations were computed using the
5 *cor()* function in the R-base *stats* package. Zero-inflated Kendall correlation and rho (ρ) were
6 computed using the *dismay()* function in the *dismay* R package⁴².

7 To make our benchmarking comparable, we adapted our clustering pipeline to remove all
8 non-automated steps and always generate a fixed number of clusters. First, hierarchical
9 clustering was performed on each correlation matrix using *dynamicTreeCut()*⁵⁵ and the
10 following non-default parameters to maximize granularity: *deepSplit = 4*, *pamStage = FALSE*,
11 *minClusterSize = 10*. Of note, we skipped the quality filtering step based on intra-cluster
12 correlations (see isoform clustering section above) to avoid bias against metrics that tend to
13 yield low values when applied to single-cell data. Since we intended to evaluate the number of
14 features remaining unclustered using each metric, we additionally suppressed the
15 unclustered isoform assignment step (see isoform clustering section above). Finally, the
16 merge process was automated by using the traditional hierarchical clustering algorithm
17 (implemented in the *hclust()* function in the R base package) to group clusters based on the
18 inferred metatranscripts that summarize the cluster's expression profile (see isoform
19 clustering section above). Finally, we set the number of clusters to 15, i.e. the number of
20 simulated co-expression modules.

21 In addition to the number of unclustered isoforms, we used the levels of internal correlation
22 in the empirical clusters, i.e. those obtained by *de novo* clustering of simulated synthetic

1 features, to evaluate the clustering. We did this by aggregately considering all pairwise metric
2 values for features within a cluster and measuring the percentage of metrics that are above a
3 threshold value of 0.8. To assess how well empirical clusters recapitulated the co-expression
4 simulation, we paired empirical with simulated clusters using the correlations between their
5 mean cluster profiles. Simulated clusters were therefore paired with the empirical clustering
6 showing maximum profile correlation. We next compared synthetic feature IDs assigned to
7 the obtained and empirical clusters in each pair using the Jaccard Index (JI).

8 **Differential Isoform Usage and co-Differential Isoform Usage across multiple** 9 **groups**

10 *Defining Differential Isoform Usage across multiple groups*

11 Grouping isoforms into different clusters allows detection of a number of expression patterns
12 across the multiple cell types included in single-cell data. As previously described, we filtered
13 DE isoforms to ensure that all transcripts had at least one other counterpart from the same
14 gene that was also significantly DE. Intuitively, in order for Differential Isoform Usage (DIU) to
15 occur, a gene must first have at least two DE isoforms. However, we only considered a gene to
16 be positive for DIU if (at least) two isoforms were DE and were assigned to different clusters,
17 indicating that two of the gene's isoforms show different expression patterns across groups
18 (see main text, [Figure 5A](#)). Ultimately, this can be interpreted as an indicator that isoform
19 expression regulation is cell type-dependent in that gene.

20 *Detecting co-splicing patterns across isoform clusters: co-Differential Isoform Usage*

1 We define coordinated splicing patterns as a situation where post-transcriptional regulation,
2 defined by isoform expression, can be detected independently of transcriptional regulation,
3 i.e. gene-level expression. To detect splicing coordination, we defined co-Differential Isoform
4 Usage (coDIU) as a pattern where a group of genes shows co-expression of their isoforms, but
5 no co-expression can be detected when only gene expression is considered (see main text,
6 **Figure 5B**). In the context of our pipeline, a set of potentially coDIU genes will have at least two
7 of their isoforms assigned to the same clusters, therefore showing detectable isoform-level
8 co-expression, and suggesting coordinated splicing regulation in that group of genes.
9 However, clustering allows expression pattern variability among members, and therefore
10 some isoforms might be assigned to clusters that do not faithfully represent their expression
11 profile, leading to detection of false-positive coDIU genes.

12 To identify groups of genes candidates for coDIU, we applied negative-binomial generalized
13 linear regression models. Let G be a group of genes, each of them with l_g isoforms, where $g =$
14 $1, \dots, |G|$. At least one of the isoforms of each gene g in G must belong to the same cluster c ,
15 where $c \in \{1, \dots, C\}$ and C is the total number of clusters. Let \mathbf{z} be the expression vector obtained
16 after concatenating the expression vectors \mathbf{y}_i of each isoform i of every gene $g = 1, \dots, |G|$. For
17 the sake of simplicity, let us assume that $|G|=2$, $l_g=2 \forall g$, and consequently $C=2$. In this case,
18 vector \mathbf{z} will contain $4N$ elements, where N is the total number of cells in the data ($N=241$ in
19 our data) and will be the response variable in our regression model. We need to assess if \mathbf{z}
20 values follow the trend depicted in **Figure 5B**, that is, the average profile across cell types of
21 the two isoforms in cluster 1 must be significantly different to the average profile of the two
22 isoforms in cluster 2. In addition, the average profile of the two isoforms of gene 1 must not
23 be different to the average profile of the two isoforms of gene 2. To identify groups of genes

1 with these characteristics we proposed to fit the regression model in Equation 4 and select
2 the group of tested genes as coDIU candidates when having a significant interaction between
3 cluster and cell type effects, and a non-significant interaction between gene and cell type
4 effects.

$$5 \quad h(z) = \beta_0 + \beta_1 G_2 + \beta_2 C_2 + \sum_{k=2}^K \gamma_k T_k + \beta_3 G_2 C_2 + \sum_{k=2}^K \delta_k T_k G_2 + \sum_{k=2}^K \tau_k T_k C_2 + \varepsilon \quad [\text{Equation 4}]$$

6 Where G_2 and C_2 are dummy variables indicating whether the expression value corresponds to
7 gene or cluster 2 (value 1) or 1 (value 0), respectively, T_k is a dummy variable which takes
8 value 1 when the corresponding cell is assigned to cell type k ($k=1, \dots, K$) and 0 otherwise, β_k , γ_k , δ_k
9 and τ_k are the regression coefficients, ε represents the error term, and $h()$ is the link function
10 of the GLM (natural logarithm in this case).

11 We fitted the GLM model with the `glm()` function in the R-base package, and the
12 `negative.binomial()` function in the *MASS* R package⁹⁴, with $\theta = 10$. To test the significance level
13 of the *cluster*cell type* and *gene*cell type* interactions, we calculated type-II analysis-of-variance
14 (ANOVA) tables for the model using a likelihood-ratio X^2 test, implemented in the *Anova()*
15 function from the *car* R package⁹⁵, since we had an unbalanced design. P -values for each of
16 the interactions were separately adjusted using the Benjamini & Hochberg correction. Gene
17 pairs were considered positive for coDIU if FDR adjusted p -value < 0.05 for the *cluster*cell type*
18 interaction and FDR adjusted p -value > 0.05 for the *gene*cell type* interaction. In other words,
19 we required expression variance across cell types to be a function of the expression profile
20 captured by the clustering, while imposing the additional limitation that aggregating
21 expression by gene must make this effect undetectable. Given that all genes with clustered

1 isoforms will form pairs with all potentially coDIU counterparts and be repeatedly tested, we
2 considered genes to be positive for coDIU if they met the significance criteria in at least one of
3 these pairwise tests.

4 **Functional analyses**

5 The analyses in this manuscript are based on a long read-defined transcriptome which, after
6 careful quality control and curation of the isoform models, was further annotated using
7 IsoAnnotLite (<https://isoannot.tappas.org/isoannot-lite/>) to include positionally-defined
8 functional features in the annotation (see [Supplementary Note 1](#)). *Functional features* are
9 grouped in *functional categories* depending on the database from which the information was
10 retrieved and on the biological functions performed by the features (comprehensive list in
11 [Supplementary Note 1](#)). In this manner, we gathered sufficient information to couple our co-
12 expression analyses with a biological readout. The specific analysis strategies used to this end
13 are detailed below.

14 *Functional Enrichment Analysis*

15 In order to understand the functional properties of AS-regulated and co-regulated genes, we
16 set out to characterize DIU and coDIU genes using different functional enrichment analysis
17 approaches. In this manner, we intended to gain insight on functional features and categories
18 showing significant overrepresentation in each of these two gene lists, in comparison to
19 different backgrounds, i.e. lists of genes to compare to in order to detect enrichment.

1 In the case of DIU genes, we calculated enrichment relative to genes with multiple DE
2 isoforms (total: 10,100) in order to discriminate the functional properties of genes regulated
3 by alternative splicing, as opposed to those lacking differential usage of their isoforms. We
4 considered all annotated functional categories and features, and applied tappAS Functional
5 Enrichment Analysis (FEA), which relies on the *GOSeg* R package⁹⁶. Briefly, the method
6 performs an over-representation Fisher's Exact test for each functional feature, considering
7 the number of genes annotated with the feature in the tests and background lists. tappAS
8 next corrects for multiple testing within each functional category by the Benjamini-Hochberg
9 method, allowing multiple functional databases to be included or excluded from the analysis
10 without influencing the number of significant features after p-value adjustment. Significant
11 enrichment for the different tests was defined using a threshold of $FDR < 0.05$.

12 For coDIU genes, we designed a different strategy in order to improve the statistical power of
13 our functional enrichment analysis, aiming to compare functional properties between splicing
14 regulation (DIU) and co-regulation (coDIU). As stated above, DIU regulation is best measured
15 by using genes with DE isoforms as background. Intuitively, coDIU-regulated genes should
16 then be characterized by comparing them to DIU genes. To accommodate these two
17 test/background lists in a functional enrichment analysis without ignoring the overlap
18 between the coDIU and DIU gene groups, we computed enrichment using a partially
19 overlapping samples z-test via the *Prop.test()* function in the *Partiallyoverlapping* R package⁶¹.
20 Specifically, we compared the proportion of coDIU genes containing each of the functional
21 features (relative to DIU genes) with the proportion of DIU genes containing that same
22 annotation (with respect to genes with DE isoforms). In other words, we tested whether the
23 proportion of coDIU vs DIU genes including a given functional feature was significantly higher

1 than that shown in the comparison between DIU and DE genes. We performed the analysis
2 for features with more than 20 annotated genes, and subsequently corrected for multiple
3 testing within functional categories using the Benjamini-Hochberg method. Functional
4 features were considered to be present in a significantly higher proportion in coDIU genes
5 when $FDR < 0.05$.

6 *Functional Diversity Analysis*

7 To obtain insight into the functional changes generated as a consequence of DIU and coDIU,
8 we again used the tappAS tool for the functional analysis of alternative splicing⁶⁰. In particular,
9 we first applied tappAS' Functional Diversity Analysis (FDA) module (see main text [Figure 6A](#)).
10 Briefly, FDA performs a within-gene comparison of all the isoforms included in the analysis,
11 aiming to detect whether they present variation in the inclusion of a functional feature. In
12 FDA, variation can be positional, i.e. one or more of the gene's isoforms present a change in
13 the genomic coordinates defining the feature, or be defined by presence/absence, i.e. at least
14 one of the isoforms lacks a feature that is present in the rest. As a result, FDA provides
15 analyzed genes with a label for each of the feature categories included in the transcriptome's
16 functional annotation file, flagging them as *varying* if at least one of the isoforms presents
17 variation in a feature from that category, or *not varying* if no changes are detected. For more
18 details on FDA, see the Methods section in de la Fuente *et al.*⁶⁰.

19 We run both positional and presence/absence FDA for three gene sets: 1) genes with multiple
20 DE isoforms (total: 3,172), 2) DIU genes (total: 2,577) and 3) coDIU genes (total: 2,049). Next,
21 for each of these gene sets, we computed the proportion of *varying* genes detected for each
22 functional category. Varying proportions were calculated relative to the total number of genes

1 including annotations from the category, instead of considering all genes in the set. In this
2 manner, we avoided underestimating variation rates for categories that were less
3 represented in the functional annotation file. In order to check whether any of these gene
4 sets presented a significantly higher mean proportion of varying genes across categories, we
5 performed a paired *t*-test for each combination of gene set pairs: DIU vs multiple DE, coDIU vs
6 multiple DE, and coDIU vs DIU. In this analysis, we considered functional categories to be the
7 individuals under evaluation, while the proportion of varying genes calculated for each
8 category in the two tested sets constituted the paired observations. As a result, we obtained
9 three p-values per FDA analysis type, i.e. presence/absence and positional variation.

10 To better understand the functional readout that can be obtained using the *acorde* pipeline,
11 we analyzed a subset of the coDIU gene network, namely three clusters showing related
12 isoform co-expression patterns: neuron-specific expression (cluster 8), oligodendrocyte-
13 specific expression (cluster 13) and expression in both neural and oligodendrocyte cell types
14 (cluster 19). In total, 160 coDIU genes showed co-expressed isoforms between at least two of
15 these clusters. To characterize functional variation among the clusters, we used
16 positional/presence FDA (see above) and ID-level FDA. ID-level FDA is also included in tappAS⁶⁰
17 and provides a within-feature summary of FDA results. In other words, ID-level FDA ultimately
18 reports the number of *varying* and *not varying* genes detected for each feature ID included in a
19 given functional category. In this case, *varying* status obeys a similar criterion to the one
20 described above, i.e. genes in which at least one isoform shows differential
21 inclusion/exclusion of the feature. Since each functional category may include several
22 features, ID-level FDA provides a complementary view to that of FDA, allowing users to inspect
23 which particular features are more frequently changing as a result of the category-level

- 1 functional variation reported in FDA. For more details on ID-level FDA, see the Methods
- 2 section in de la Fuente *et al.*⁶⁰.

1 References

- 2 1. Svensson, V., Vento-Tormo, R. & Teichmann, S. A. Exponential scaling of single-cell RNA-seq in the
3 past decade. *Nat. Protoc.* **13**, 599–604 (2018).
- 4 2. Treutlein, B. *et al.* Reconstructing lineage hierarchies of the distal lung epithelium using single-cell
5 RNA-seq. *Nature* **509**, 371–375 (2014).
- 6 3. Zeisel, a. *et al.* Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq.
7 *Science (80-.)*. **347**, 1138–42 (2015).
- 8 4. Chen, R., Wu, X., Jiang, L. & Zhang, Y. Single-Cell RNA-Seq Reveals Hypothalamic Cell Diversity. *Cell*
9 *Rep.* **18**, 3227–3241 (2017).
- 10 5. Wu, Y. E., Pan, L., Zuo, Y., Li, X. & Hong, W. Detecting Activated Cell Populations Using Single-Cell
11 RNA-Seq. *Neuron* **96**, 313–329.e6 (2017).
- 12 6. Villani, A.-C. *et al.* Single-cell RNA-seq reveals new types of human blood dendritic cells,
13 monocytes, and progenitors. *Science (80-.)*. **356**, eaah4573 (2017).
- 14 7. Zhong, S. *et al.* A single-cell RNA-seq survey of the developmental landscape of the human
15 prefrontal cortex. *Nature* **555**, 524–528 (2018).
- 16 8. Crow, M. & Gillis, J. Co-expression in Single-Cell Analysis: Saving Grace or Original Sin? *Trends in*
17 *Genetics* vol. 34 823–831 (2018).
- 18 9. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
19 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
- 20 10. Saelens, W., Cannoodt, R., Todorov, H. & Saeys, Y. A comparison of single-cell trajectory inference
21 methods. *Nat. Biotechnol.* **37**, 547–554 (2019).
- 22 11. La Manno, G. *et al.* RNA velocity of single cells. *Nature* **560**, 494–498 (2018).
- 23 12. Jia, G. *et al.* Single cell RNA-seq and ATAC-seq analysis of cardiac progenitor cell transition states
24 and lineage settlement. *Nat. Commun.* **9**, 4877 (2018).
- 25 13. Su, X. *et al.* Single-cell RNA-Seq analysis reveals dynamic trajectories during mouse liver
26 development. *BMC Genomics* **18**, 946 (2017).
- 27 14. Guo, F. *et al.* Single-cell multi-omics sequencing of mouse early embryos and embryonic stem cells.
28 *Cell Res.* **27**, 967–988 (2017).
- 29 15. Le, J. *et al.* Single-Cell RNA-Seq Mapping of Human Thymopoiesis Reveals Lineage Specification
30 Trajectories and a Commitment Spectrum in T Cell Development. *Immunity* **52**, 1105–1118.e9
31 (2020).
- 32 16. Jerber, J. *et al.* Population-scale single-cell RNA-seq profiling across dopaminergic neuron
33 differentiation. *Nat. Genet.* **53**, 304–312 (2021).
- 34 17. Pratapa, A., Jalihal, A. P., Law, J. N., Bharadwaj, A. & Murali, T. M. Benchmarking algorithms for gene
35 regulatory network inference from single-cell transcriptomic data. *Nat. Methods* **17**, 147–154
36 (2020).
- 37 18. Westoby, J., Artemov, P., Hemberg, M. & Ferguson-Smith, A. Obstacles to detecting isoforms using
38 full-length scRNA-seq data. *Genome Biol.* **21**, 74 (2020).
- 39 19. Arzalluz-Luque, Á. & Conesa, A. Single-cell RNAseq for the study of isoforms—how is that possible?
40 *Genome Biol.* **19**, 110 (2018).
- 41 20. Ntranos, V., Yi, L., Melsted, P. & Pachter, L. A discriminative learning approach to differential
42 expression analysis for single-cell RNA-seq. *Nat. Methods* **16**, 163–166 (2019).

- 1 21. Qiu, X. *et al.* Single-cell mRNA quantification and differential analysis with Census. *Nat. Methods*
2 **14**, 309–315 (2017).
- 3 22. Song, Y. *et al.* Single-Cell Alternative Splicing Analysis with Expedition Reveals Splicing Dynamics
4 during Neuron Differentiation. *Mol. Cell* **67**, 148–161.e5 (2017).
- 5 23. Huang, Y. & Sanguinetti, G. BRIE: transcriptome-wide splicing quantification in single cells. *Genome*
6 *Biol.* **18**, 123 (2017).
- 7 24. Wu, X., Liu, T., Ye, C., Ye, W. & Ji, G. scAPAtrop: identification and quantification of alternative
8 polyadenylation sites from single-cell RNA-seq data. *Brief. Bioinform.* **2020**, 1–15 (2020).
- 9 25. Hu, Y., Wang, K. & Li, M. Detecting differential alternative splicing events in scRNA-seq with or
10 without Unique Molecular Identifiers. *PLoS Comput. Biol.* **16**, e1007925 (2020).
- 11 26. Patrick, R. *et al.* Sierra: Discovery of differential transcript usage from polyA-captured single-cell
12 RNA-seq data. *Genome Biol.* **21**, 1–27 (2020).
- 13 27. Boeshaghi, A. S. *et al.* Isoform cell type specificity in the mouse primary motor cortex. *bioRxiv*
14 (2020) doi:10.1101/2020.03.05.977991.
- 15 28. Byrne, A. *et al.* Nanopore long-read RNAseq reveals widespread transcriptional variation among the
16 surface receptors of individual B cells. *Nat. Commun.* **8**, 16027 (2017).
- 17 29. Volden, R. *et al.* Improving nanopore read accuracy with the R2C2 method enables the sequencing
18 of highly multiplexed full-length single-cell cDNA. *Proc. Natl. Acad. Sci.* **115**, 9726–9731 (2018).
- 19 30. Gupta, I. *et al.* Single-cell isoform RNA sequencing characterizes isoforms in thousands of
20 cerebellar cells. *Nat. Biotechnol.* **36**, 1197–1202 (2018).
- 21 31. Joglekar, A. *et al.* A spatially resolved brain region- and cell type-specific isoform atlas of the
22 postnatal mouse brain. *Nat. Commun.* **12**, 463 (2021).
- 23 32. Tian, L. *et al.* Comprehensive characterization of single cell full-length isoforms in human and
24 mouse with long-read sequencing. *bioRxiv* (2020) doi:10.1101/2020.08.10.243543.
- 25 33. Feng, H. *et al.* Complexity and graded regulation of neuronal cell-type-specific alternative splicing
26 revealed by single-cell RNA sequencing. *Proc. Natl. Acad. Sci.* **118**, e2013056118 (2021).
- 27 34. Becht, E., Zhao, E., Amezcua, R. & Gottardo, R. Aggregating transcript-level analyses for single-
28 cell differential gene expression. *Nat. Methods* **17**, 583–585 (2020).
- 29 35. Yi, L., Pimentel, H., Bray, N. L. & Pachter, L. Gene-level differential analysis at transcript-level
30 resolution. *Genome Biol.* **19**, 53 (2018).
- 31 36. Soneson, C. & Robinson, M. D. Bias, robustness and scalability in single-cell differential expression
32 analysis. *Nat. Methods* **15**, 255–261 (2018).
- 33 37. Shalek, A. K. *et al.* Single-cell transcriptomics reveals bimodality in expression and splicing in
34 immune cells. *Nature* **498**, 236–240 (2013).
- 35 38. Liu, W. & Zhang, X. Single-cell alternative splicing analysis reveals dominance of single transcript
36 variant. *Genomics* **112**, 2418–2425 (2020).
- 37 39. Buen Abad Najar, C. F., Yosef, N. & Lareau, L. F. Coverage-dependent bias creates the appearance
38 of binary splicing in single cells. *Elife* **9**, 1–23 (2020).
- 39 40. Tilgner, H. *et al.* Microfluidic isoform sequencing shows widespread splicing coordination in the
40 human transcriptome. *Genome Res.* **28**, 231–242 (2018).
- 41 41. Tilgner, H. *et al.* Comprehensive transcriptome analysis using synthetic long-read sequencing
42 reveals molecular co-association of distant splicing events. *Nat. Biotechnol.* **33**, 736–742 (2015).
- 43 42. Skinnider, M. A., Squair, J. W. & Foster, L. J. Evaluating measures of association for single-cell
44 transcriptomics. *Nat. Methods* **16**, 381–386 (2019).

- 1 43. Tasic, B. *et al.* Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nat.*
2 *Neurosci.* **19**, 335–346 (2016).
- 3 44. Wyman, D. *et al.* A technology-agnostic long-read analysis pipeline for transcriptome discovery and
4 quantification. *bioRxiv* (2019) doi:10.1101/672931.
- 5 45. Soneson, C., Matthes, K. L., Nowicka, M., Law, C. W. & Robinson, M. D. Isoform prefiltering improves
6 performance of count-based methods for analysis of differential transcript usage. *Genome Biol.* **17**,
7 12 (2016).
- 8 46. Tardaguila, M. *et al.* SQANTI: extensive characterization of long-read transcript sequences for
9 quality control in full-length transcriptome identification and quantification. *Genome Res.* **28**, 396–
10 411 (2018).
- 11 47. Van den Berge, K. *et al.* Observation weights unlock bulk RNA-seq tools for zero inflation and
12 single-cell applications. *Genome Biol.* **19**, 24 (2018).
- 13 48. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq
14 data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 15 49. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: A Bioconductor package for differential
16 expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140 (2009).
- 17 50. Chen, S. & Mar, J. C. Evaluating methods of inferring gene regulatory networks highlights their lack
18 of performance for single cell gene expression data. *BMC Bioinformatics* **19**, 232 (2018).
- 19 51. Iacono, G., Massoni-Badosa, R. & Heyn, H. Single-cell transcriptomics unveils gene regulatory
20 network plasticity. *Genome Biol.* **20**, 110 (2019).
- 21 52. Brennecke, P. *et al.* Accounting for technical noise in single-cell RNA-seq experiments. *Nat. Methods*
22 **10**, 1093–1095 (2013).
- 23 53. Grün, D., Kester, L. & van Oudenaarden, A. Validation of noise models for single-cell
24 transcriptomics. *Nat. Methods* **11**, 637–640 (2014).
- 25 54. Raj, A. & van Oudenaarden, A. Nature, Nurture, or Chance: Stochastic Gene Expression and Its
26 Consequences. *Cell* vol. 135 216–226 (2008).
- 27 55. Langfelder, P., Zhang, B. & Horvath, S. Defining clusters from a hierarchical cluster tree: The
28 Dynamic Tree Cut package for R. *Bioinformatics* **24**, 719–720 (2008).
- 29 56. Pimentel, R. S., Niewiadomska-Bugaj, M. & Wang, J.-C. Association of zero-inflated continuous
30 variables. *Stat. Probab. Lett.* **96**, 61–67 (2015).
- 31 57. Erb, I. & Notredame, C. How should we measure proportionality on relative gene expression data?
32 *Theory Biosci.* **135**, 21–36 (2016).
- 33 58. Zhang, X., Xu, C. & Yosef, N. Simulating multiple faceted variability in single cell RNA sequencing.
34 *Nat. Commun.* **10**, 2611 (2019).
- 35 59. Vitting-Seerup, K. & Sandelin, A. IsoformSwitchAnalyzeR: analysis of changes in genome-wide
36 patterns of alternative splicing and its functional consequences. *Bioinformatics* **35**, 4469–4471
37 (2019).
- 38 60. de la Fuente, L. *et al.* tappAS: a comprehensive computational framework for the analysis of the
39 functional impact of differential splicing. *Genome Biol.* **21**, 119 (2020).
- 40 61. Derrick, B., White, P. & Toher, D. Parametric and non-parametric tests for the comparison of two
41 samples which both include paired and unpaired observations. *J. Mod. Appl. Stat. Methods* **18**, 2–23
42 (2019).
- 43 62. Ni, J. Z. *et al.* Ultraconserved elements are associated with homeostatic control of splicing
44 regulators by alternative splicing and nonsense-mediated decay. *Genes Dev.* **21**, 708–718 (2007).

- 1 63. Damianov, A. & Black, D. L. Autoregulation of Fox protein expression to produce dominant negative
2 splicing factors. *RNA* **16**, 405–416 (2010).
- 3 64. Jangi, M., Boutz, P. L., Paul, P. & Sharp, P. A. Rbfox2 controls autoregulation in RNA-binding protein
4 networks. *Genes Dev.* **28**, 637–651 (2014).
- 5 65. Kino, Y. *et al.* Nuclear localization of MBNL1: Splicing-mediated autoregulation and repression of
6 repeat-derived aberrant proteins. *Hum. Mol. Genet.* **24**, 740–756 (2015).
- 7 66. Boutla, A., Delidakis, C. & Tabler, M. Developmental defects by antisense-mediated inactivation of
8 micro-RNAs 2 and 13 in Drosophila and the identification of putative target genes. *Nucleic Acids
9 Res.* **31**, 4973–4980 (2003).
- 10 67. Lai, E. C., Burks, C. & Posakony, J. W. The K box, a conserved 3' UTR sequence motif, negatively
11 regulates accumulation of Enhancer of split Complex transcripts. *Development* **125**, 4077–4088
12 (1998).
- 13 68. Waragai, M. *et al.* Ataxin 10 induces neurogenesis via interaction with G-protein $\beta 2$ subunit. *J.
14 Neurosci. Res.* **83**, 1170–1178 (2006).
- 15 69. Westbrook, T. F. *et al.* SCF β -TRCP controls oncogenic transformation and neural differentiation
16 through REST degradation. *Nature* **452**, 370–374 (2008).
- 17 70. Raynaud-Messina, B. & Merdes, A. γ -tubulin complexes and microtubule organization. *Current
18 Opinion in Cell Biology* vol. 19 24–30 (2007).
- 19 71. Dráberová, E. *et al.* Differential expression of human γ -tubulin isoforms during neuronal
20 development and oxidative stress points to a γ -tubulin-2 prosurvival function. *FASEB J.* **31**, 1828–
21 1846 (2017).
- 22 72. Hausrat, T. J., Radwitz, J., Lombino, F. L., Breiden, P. & Kneussel, M. Alpha- and beta-tubulin isoforms
23 are differentially expressed during brain development. *Dev. Neurobiol.* **81**, 333–350 (2020).
- 24 73. Chakraborti, S., Natarajan, K., Curiel, J., Janke, C. & Liu, J. The emerging role of the tubulin code:
25 From the tubulin molecule to neuronal function and disease. *Cytoskeleton* vol. 73 521–550 (2016).
- 26 74. Janke, C. & Magiera, M. M. The tubulin code and its role in controlling microtubule properties and
27 functions. *Nature Reviews Molecular Cell Biology* vol. 21 307–326 (2020).
- 28 75. Verhey, K. J. & Gaertig, J. The tubulin code. *Cell Cycle* vol. 6 2152–2160 (2007).
- 29 76. Tischfield, M. A. & Engle, E. C. Distinct α - and β -tubulin isoforms are required for the positioning,
30 differentiation and survival of neurons: New support for the 'multi-tubulin' hypothesis. *Bioscience
31 Reports* vol. 30 319–330 (2010).
- 32 77. Park, J. H. & Roll-Mecak, A. The tubulin code in neuronal polarity. *Current Opinion in Neurobiology*
33 vol. 51 95–102 (2018).
- 34 78. Fu, X. D. & Ares, M. Context-dependent control of alternative splicing by RNA-binding proteins.
35 *Nature Reviews Genetics* vol. 15 689–701 (2014).
- 36 79. Saha, A. *et al.* Co-expression networks reveal the tissue-specific regulation of transcription and
37 splicing. *Genome Res.* **27**, 1843–1858 (2017).
- 38 80. Aghamirzaie, D., Collakova, E., Li, S. & Grene, R. CoSpliceNet: a framework for co-splicing network
39 inference from transcriptomics data. *BMC Genomics* **17**, 845 (2016).
- 40 81. Zhang, P., Southey, B. R. & Rodriguez-Zas, S. L. Co-expression networks uncover regulation of
41 splicing and transcription markers of disease. in *EPiC Series in Computing* vol. 70 119–128 (2020).
- 42 82. Chau, K. *et al.* Isoform transcriptome of developing human brain provides new insights into autism
43 risk variants. *bioRxiv* (2020) doi:10.1101/2020.06.27.175489.

- 1 83. Vu, T. N. *et al.* Isoform-level gene expression patterns in single-cell RNA-sequencing data.
2 *Bioinforma. doi* **10**, 1-9 (2018).
- 3 84. Yap, K., Xiao, Y., Friedman, B. A., Je, H. S. & Makeyev, E. V. Polarizing the Neuron through Sustained
4 Co-expression of Alternatively Spliced Isoforms. *Cell Rep.* **15**, 1316-1328 (2016).
- 5 85. Ma, J. *et al.* Comprehensive expression-based isoform biomarkers predictive of drug responses
6 based on isoform co-expression networks and clinical data. *Genomics* **112**, 647-658 (2020).
- 7 86. Ma, J.-Q. *et al.* Differential Alternative Splicing Genes and Isoform Regulation Networks of Rapeseed
8 (*Brassica napus* L.) Infected with *Sclerotinia sclerotiorum*. *Genes (Basel)*. **11**, 784 (2020).
- 9 87. Bray, N. The power of 3' UTRs. *Nat. Rev. Neurosci.* (2018) doi:10.1038/s41583-018-0011-6.
- 10 88. Bae, B. & Miura, P. Emerging roles for 3' UTRs in neurons. *Int. J. Mol. Sci.* **21**, 3413 (2020).
- 11 89. Picelli, S. *et al.* Tn5 transposase and tagmentation procedures for massively scaled sequencing
12 projects. *Genome Res.* **24**, 2033-2040 (2014).
- 13 90. Dobin, A. *et al.* STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2013).
- 14 91. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a
15 reference genome. *BMC Bioinformatics* **12**, 323 (2011).
- 16 92. Tarazona, S. *et al.* Data quality aware analysis of differential expression in RNA-seq with NOISeq
17 R/Bioc package. *Nucleic Acids Res.* **43**, gkv711 (2015).
- 18 93. R Core Team. R: A language and environment for statistical computing. *R Foundation for Statistical
19 Computing* (2021).
- 20 94. Venables, W. & Ripley, B. *Modern Applied Statistics with S.* (Springer, 2002).
- 21 95. Fox, J. & Weisberg, S. *An R Companion to Applied Regression.* (Sage, 2019).
- 22 96. Young, M. D., Wakefield, M. J., Smyth, G. K. & Oshlack, A. Gene ontology analysis for RNA-seq:
23 accounting for selection bias. *Genome Biol.* **11**, R14 (2010).

1 **Author contributions**

2 AAL developed and evaluated computational methods, implemented and documented R
3 package, obtained long read transcriptome, analyzed single-cell data, interpreted results and
4 wrote the manuscript. PS implemented IsoAnnotLite, refined and performed functional
5 annotation strategy. ST conducted statistical method development, refined the manuscript
6 and supervised the study. AC conceived the study, supervised all analysis approaches and
7 methods developed, and refined the manuscript.

8 **Acknowledgements**

9 This work has been funded by NIH grant R21HG011280 and by the Spanish Ministry of Science
10 grants BIO2015-71658 and BES-2016-076994.

11 **Competing interests**

12 The authors declare no competing interests.

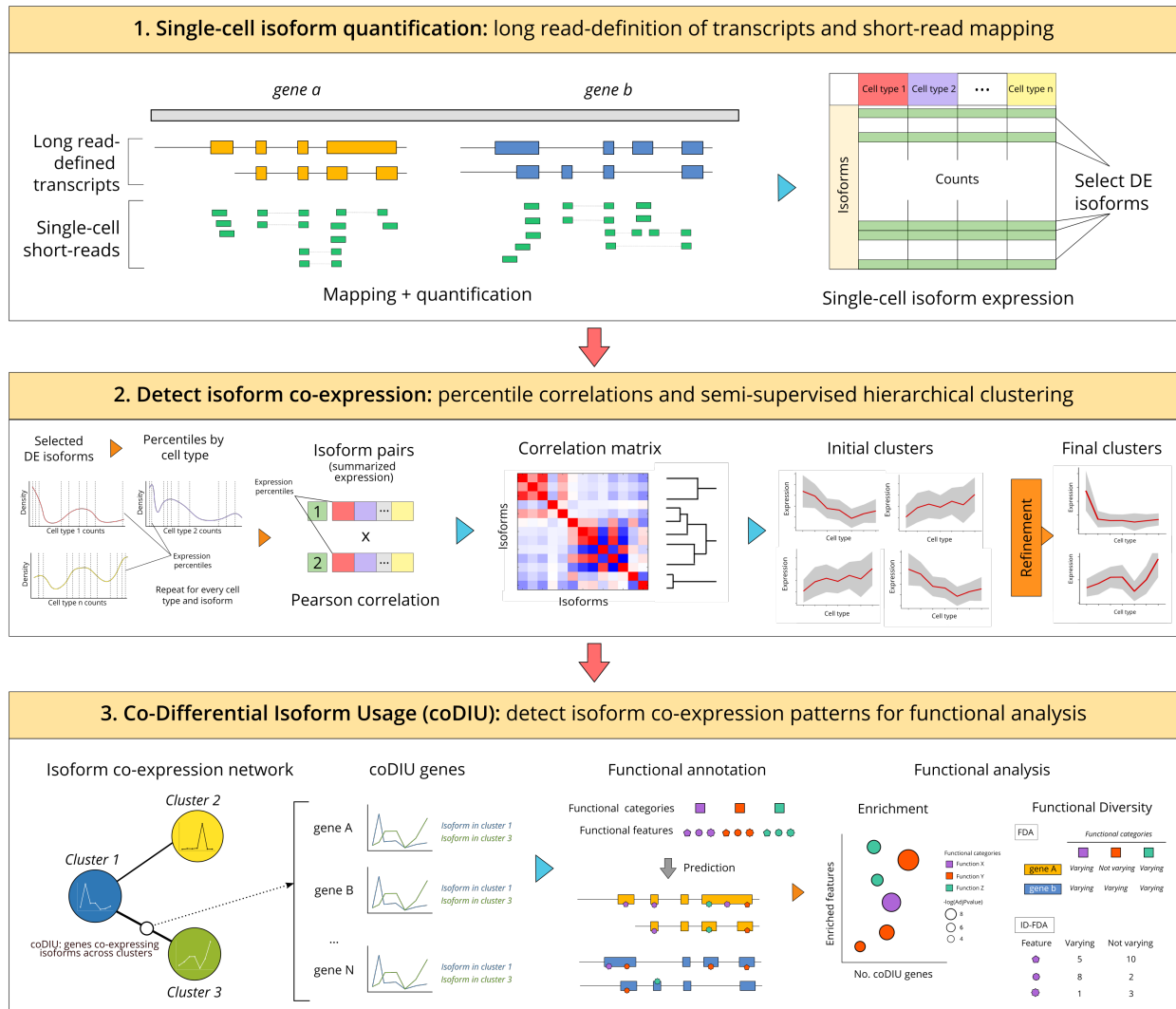


Figure 1: *acorde* workflow. The *acorde* pipeline includes three main analysis modules. First, long read RNA-Seq data is used to define isoform models and short, single-cell RNA-Seq reads are mapped to the long read-generated transcriptome. Isoform are then tested for multi-group differential expression and those that are significantly DE in at least one of the cell types are selected. Next, percentile correlations are computed to cluster isoforms with similar expression patterns across cell types. Finally, gene pairs are tested for co-differential isoform usage, detecting genes that form co-expression relationships for subsequent functional analysis.

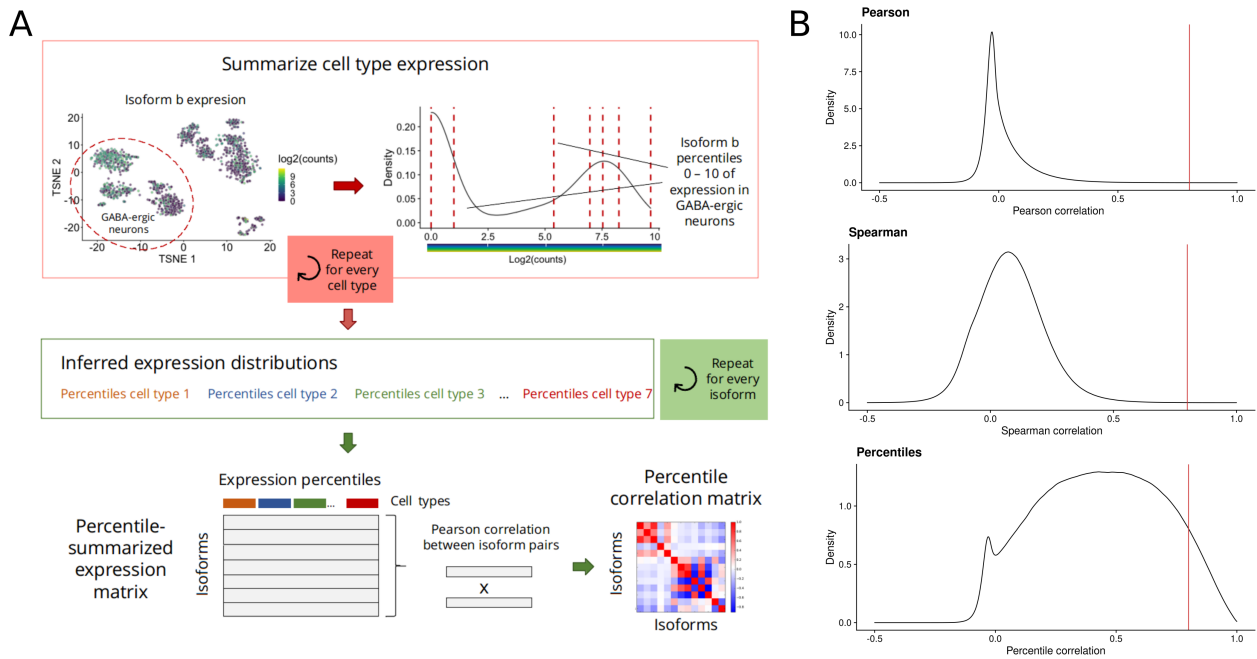


Figure 2: Percentile correlations. A) Percentile correlation algorithm. For each isoform, cell type-level expression is summarized using percentiles (0 - 10) as a proxy of the expression distribution of the isoform in each of the cell types. Then, Pearson correlations are computed using the percentile-summarized expression of all isoforms, obtaining a percentile correlation matrix. B) Correlation density distributions. Pairwise isoform correlations were computed using Pearson, Spearman and percentile+Pearson correlation. Red vertical line indicates the $cor = 0.8$ threshold.

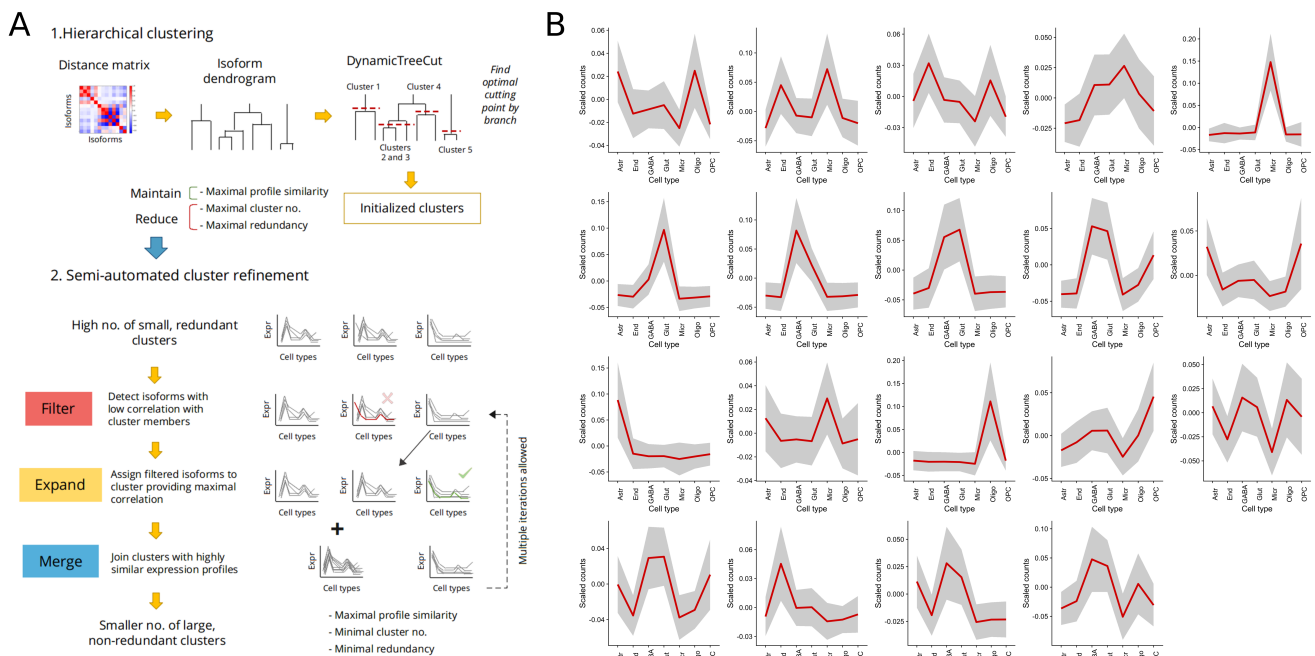


Figure 3: isoform clustering. A) Clustering pipeline. The percentile correlation matrix is first used as a distance matrix for hierarchical clustering. After dynamic cluster generation, noisy clusters are refined by a three-step semi-automated process. B) Clusters generated after applying the *acorde* clustering pipeline to the mouse neural dataset. Cell-level mean expression (scaled, see Methods) is computed for all transcripts and then aggregated as the global cell type mean, represented by the red line. Grey area corresponds to cell type mean \pm standard deviation.

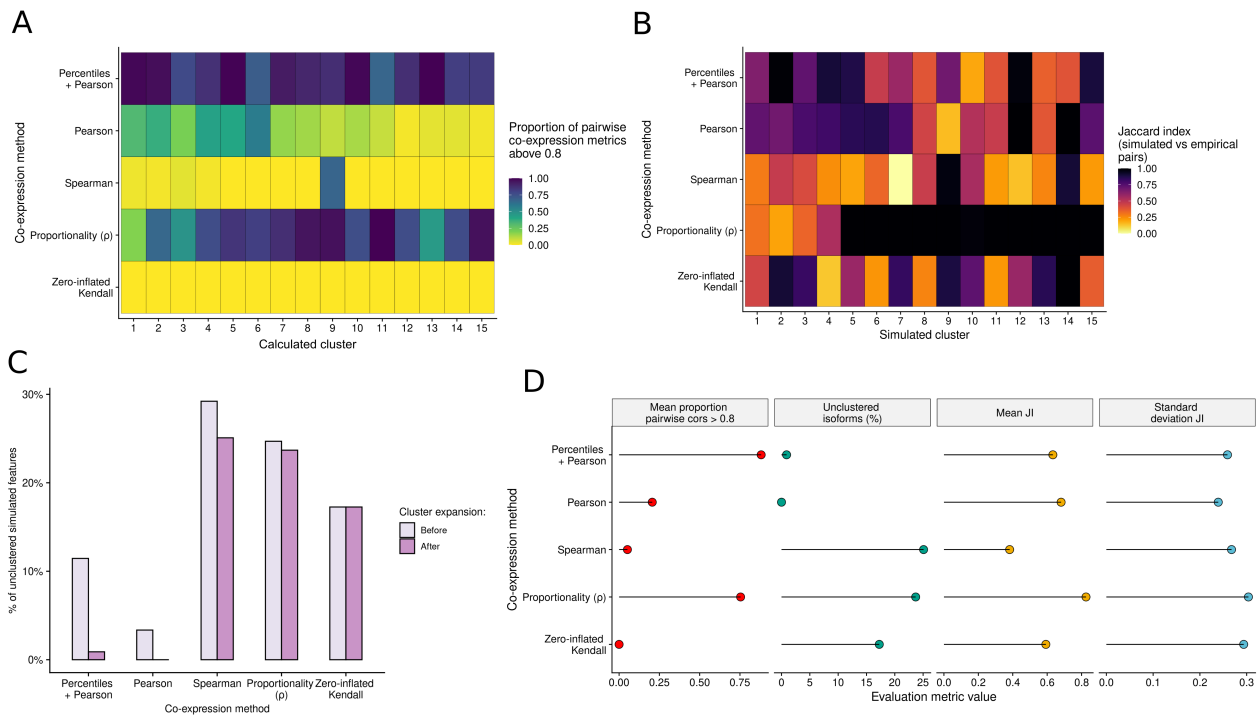


Figure 4: evaluation of percentile correlation-based clustering on simulated data. A) Proportion of co-expression values above 0.8 in each empirical cluster obtained after running the *acorde* clustering pipeline on simulated data using several correlation and proportionality methods as a distance metric. B) Jaccard Index of simulated vs calculated clusters obtained with each evaluated co-expression method. Simulated clusters were paired with one calculated cluster based on mean profile similarity, and synthetic transcripts present in each member of the pair were compared. C) Percentage of unclustered isoforms generated by each co-expression method, before and after re-assigning unclustered isoforms by measuring co-expression with the mean profile of extant clusters (i.e. cluster expansion). D) Evaluation metric overview. Metrics are specified in the grid headers. x-axis shows values of the different metrics, y-axis displays evaluated co-expression methods.

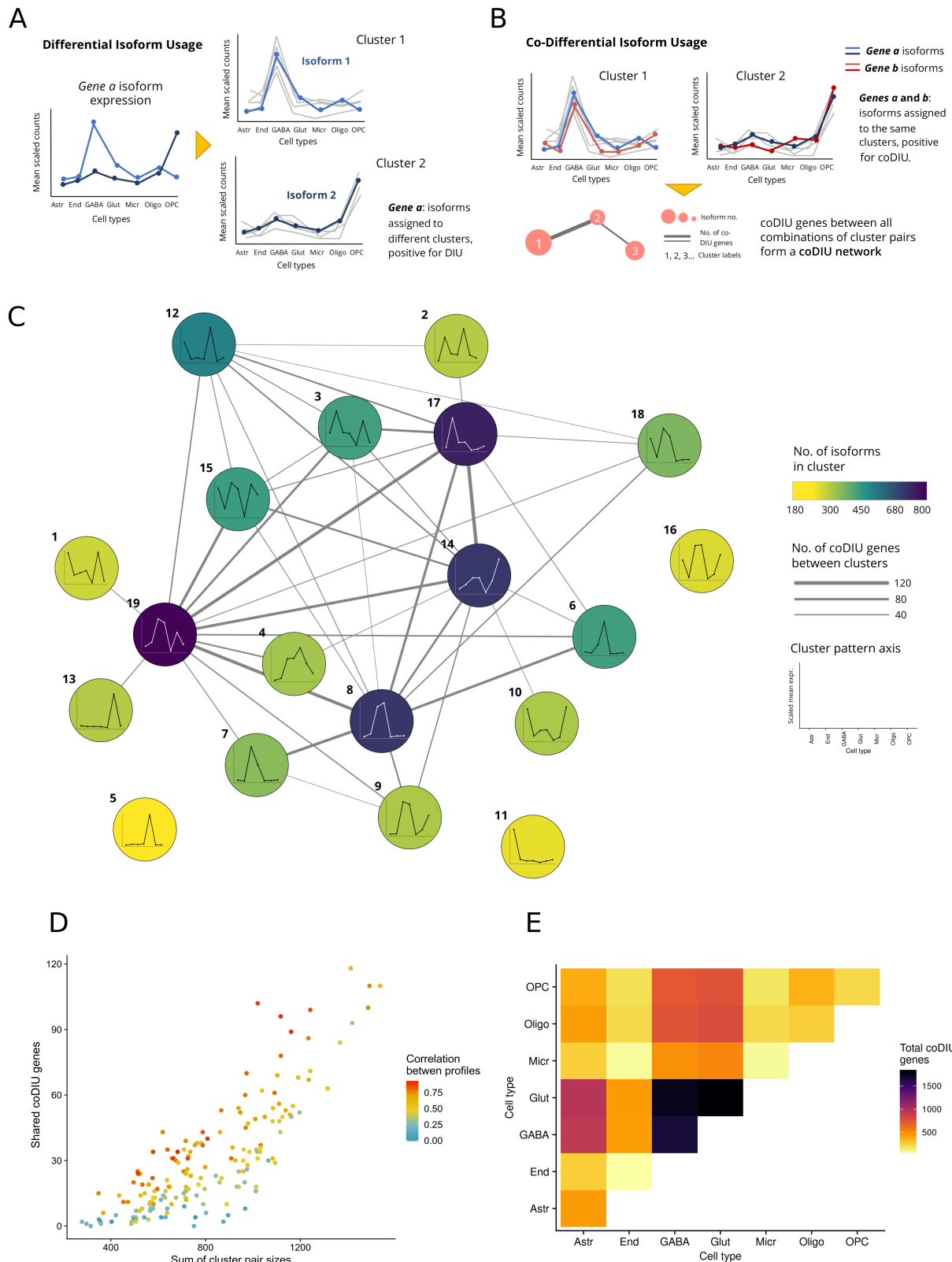


Figure 5: characterization of genes with co-Differential Isoform Usage. A) Cluster-based definition of Differential Isoform Usage (DIU) across multiple cell types. DIU genes that have at least two isoforms assigned to different clusters, indicating a differential isoform selection pattern across the different cell types. B) Definition of co-Differential Isoform Usage (coDIU) using clusters. CoDIU genes have multiple isoforms assigned to the same clusters, establishing cross-cell type co-expression relationships for at least two of their isoforms. C) coDIU network. Nodes represent clusters and depict their mean expression profile across cell types. Node color represents cluster size (i.e. no. of isoforms in cluster). Edge width represents number of coDIU genes detected between each pair of clus-

ters. D) Evaluation of cluster profile similarity and size as a function of the number of coDIU genes detected by *acorde*. x-axis corresponds to the sum of isoforms in each possible pair of clusters generated from the data. y-axis contains the number of coDIU genes between the pair. Dot color represents the correlation between the mean expression profiles of each pair of clusters. The number of coDIU genes between a pair of clusters is seemingly related to the size of the clusters involved, and shows no relationship with the degree of similarity between the expression profiles of clustered isoforms. E) Cell type-level coDIU patterns. For each pair of cell types represented in x and y-axis, heatmap color corresponds to the total number of genes found to be co-DIU between them. Total coDIU genes are calculated as the sum of coDIU genes detected between all occurrences of cluster pairs showing high expression of isoforms these cell types. GABA and Glut cell types share the highest number of coDIU genes, both with each other and with other cell types. Astr: astrcytes, End: endothelial cells, GABA: GABA-ergic neurons, Glut: glutamatergic neurons, Micr: microglia, Oligo: oligodendrocytes, OPC: oligodendrocyte precursor cells.

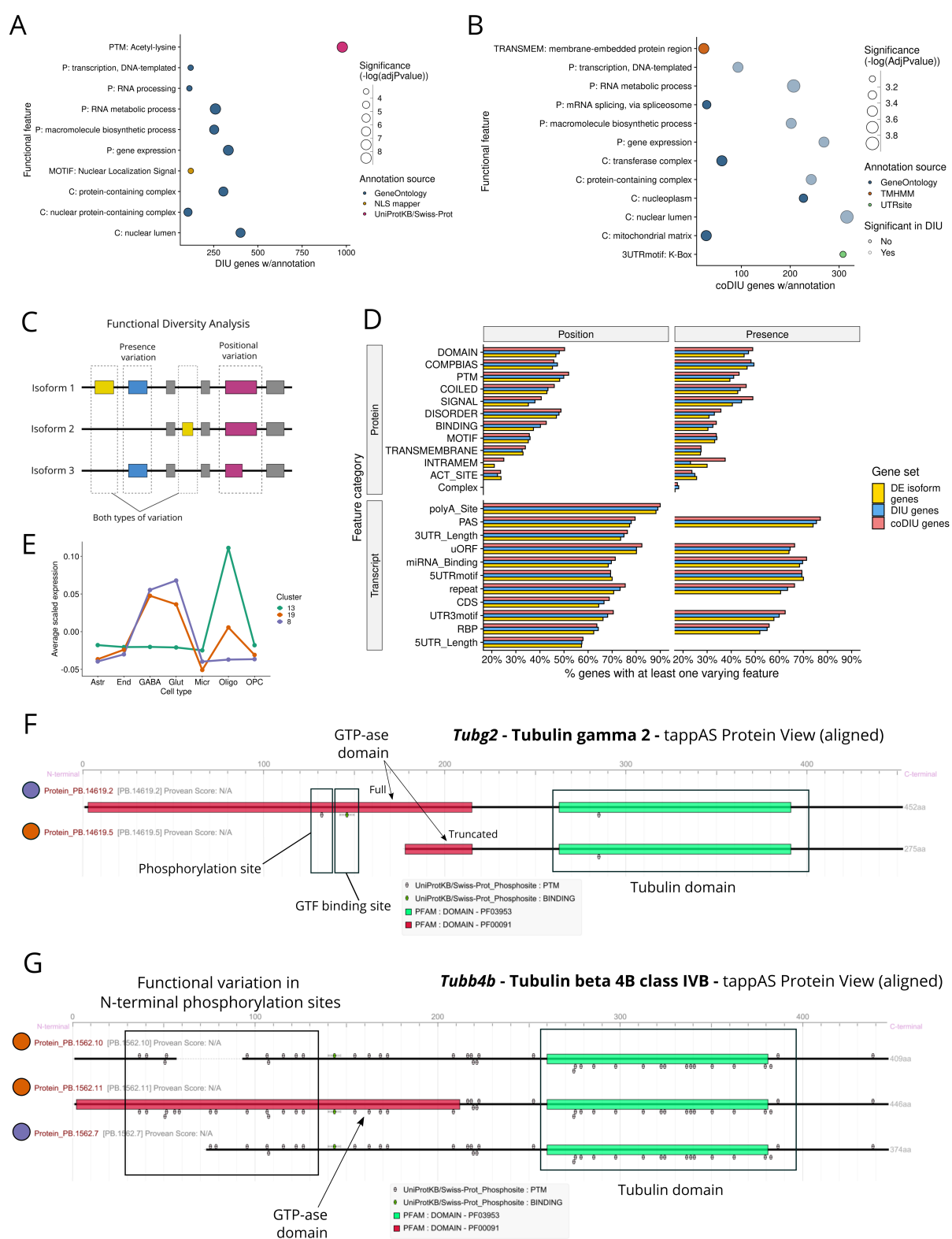
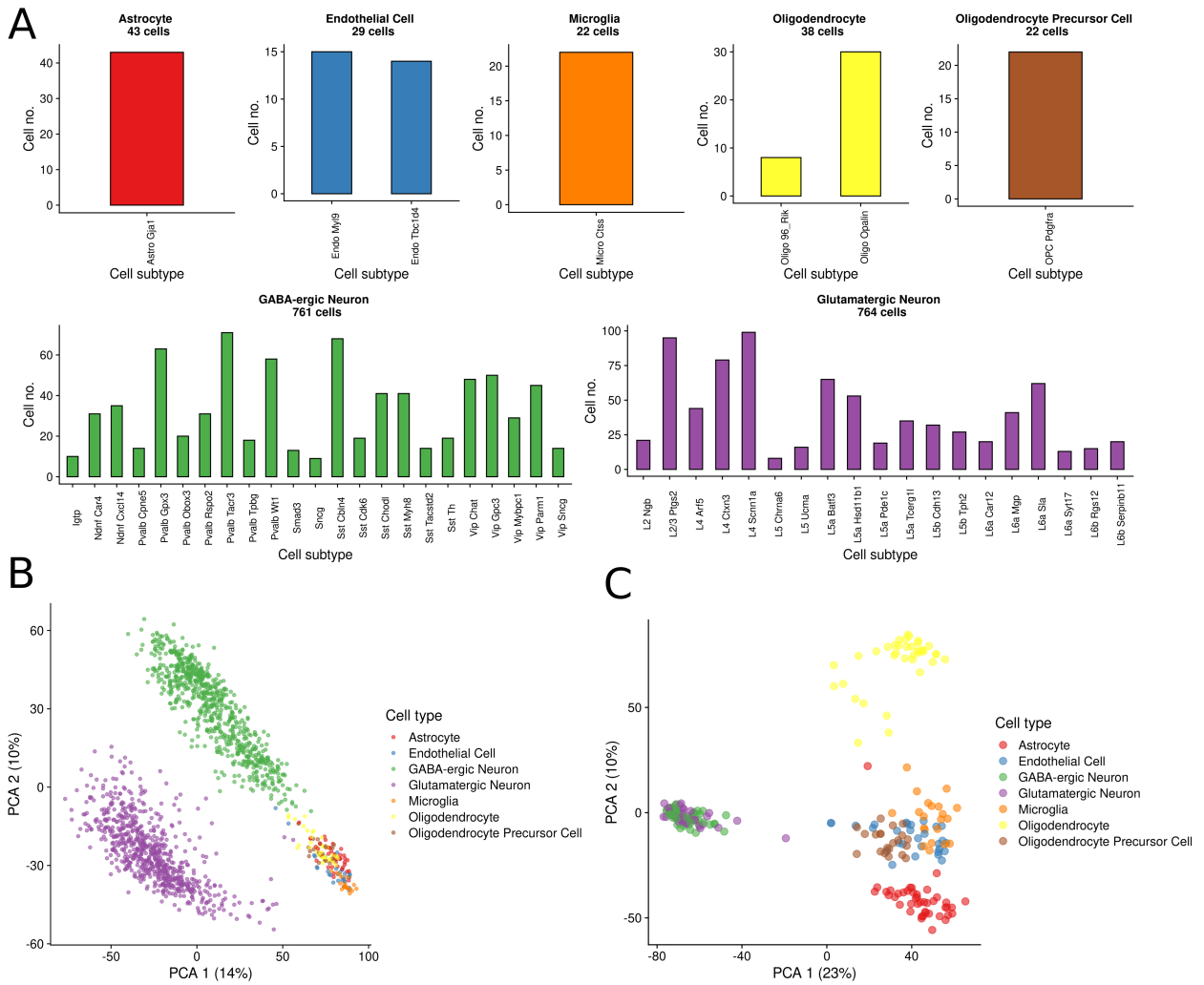
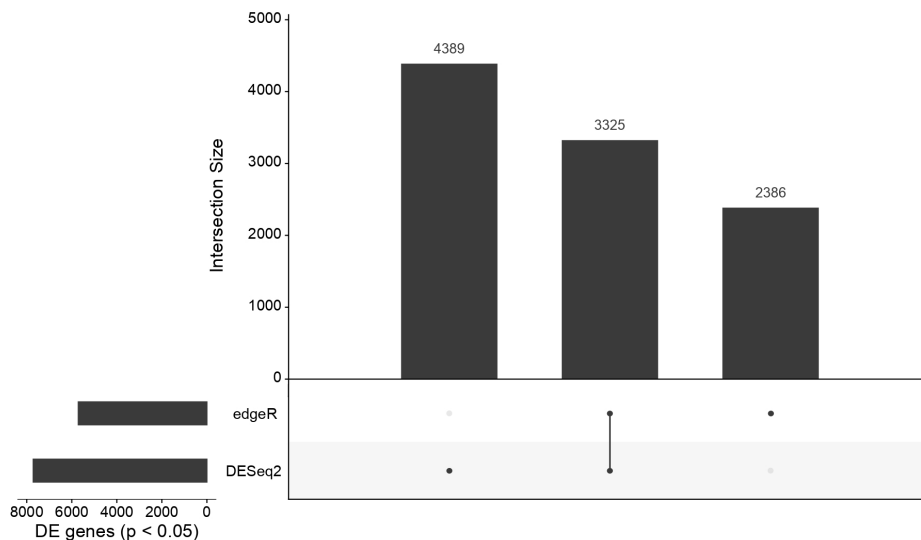


Figure 6: functional analysis of DIU and coDIU genes. Functional Enrichment results for A) Differential isoform Usage (DIU) genes vs genes with Differential Expression (DE) of isoforms and B) coDIU vs DIU genes. x-axis indicates the total number of test genes (A-DIU, B-coDIU) including the tested annotation feature, y-axis shows functional features. Dot color represents functional category, dot size represents $-\log(\text{adjusted } p\text{-value})$. C) Schematic representation of the Functional Diversity Analysis (FDA). Variation using the genomic position and present/absence criteria are shown. D) FDA results for DE isoform, DIU and coDIU genes. y-axis shows transcript and protein functional categories (see Supplementary Note 1 for category definition information). x-axis shows the percentage of genes including at least one feature annotation from each of the categories that are detected as functionally *varying*. Both FDA criteria are shown (position - left grid column, presence - right grid column). CoDIU genes show the largest le-

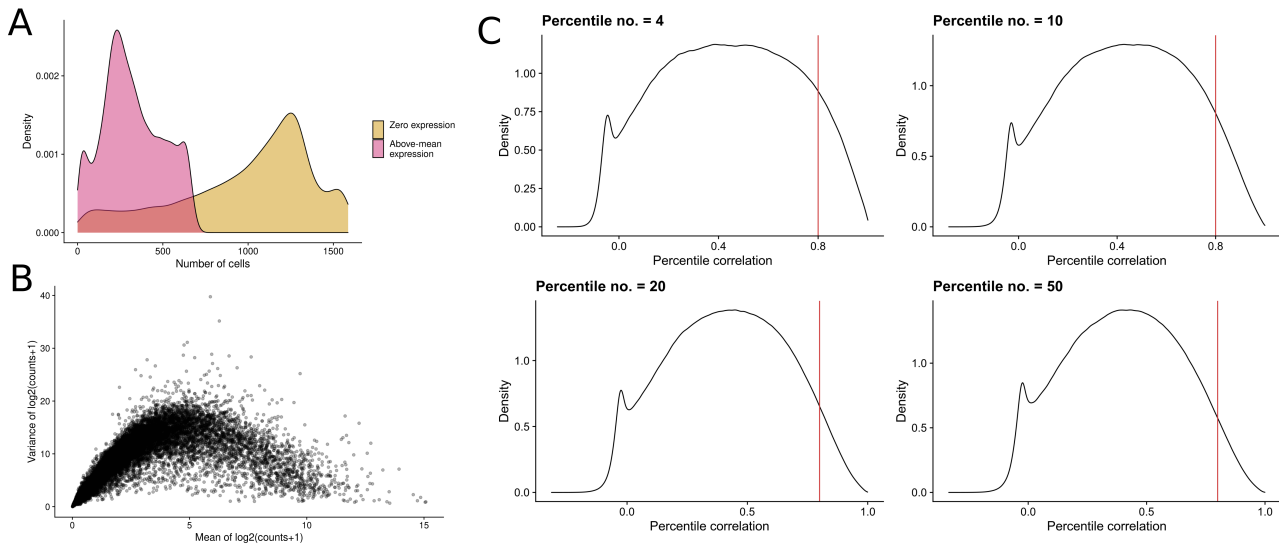
vel of functional variation. E) Cell type expression patterns of clusters selected for downstream functional analysis: neural (cluster 8, purple), oligodendrocyte (cluster 13, green) or shared (cluster 19, orange). tappAS view of F) *Tub2g* and G) *Tubb4b* protein annotations. Cluster assignments for each isoform are indicated by dot color.



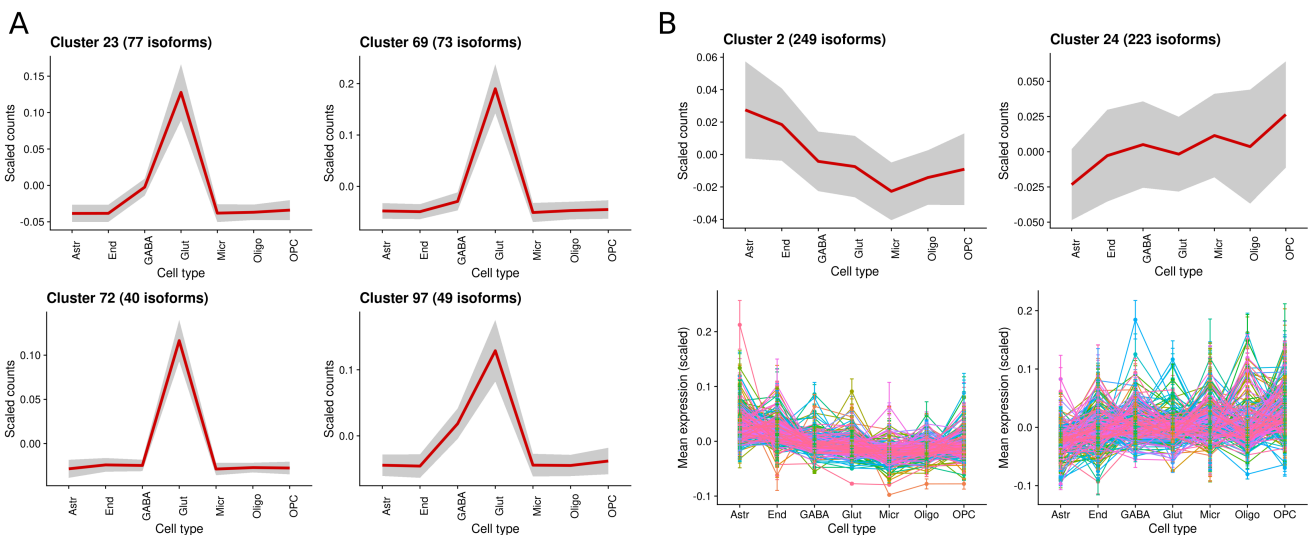
Supplementary Figure 1: Tasic et al. mouse neural dataset overview. A) Number of cells assigned to each cell type and subtype by Tasic et al. (post-QC, 1591 cells). B) Principal Component Analysis of the full dataset, components 1 and 2. C) Principal Component Analysis of data after balancing cell number by randomly sampling 50 GABA and 50 glutamatergic neural cells (242 total cells).



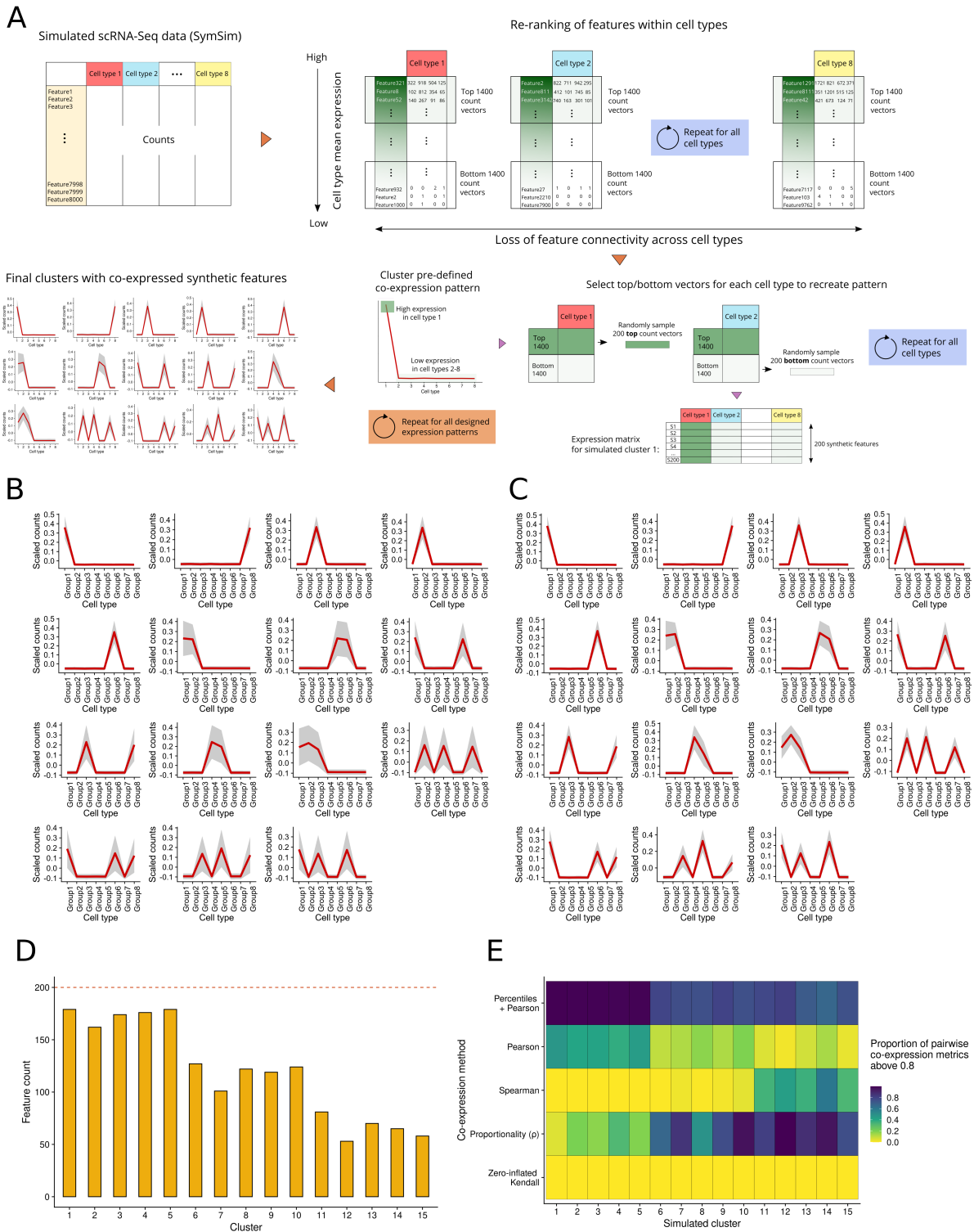
Supplementary Figure 2: upset plot of isoform-level Differential Expression (DE) analysis results. DE isoforms common and uniquely-detected by edgeR and DESeq2 are represented by vertical bar height. Analyzed sets (unique isoforms or intersection) are indicated by the dots below. Total DE isoforms detected by each method (adjusted p -value < 0.05) is represented by horizontal bar height.



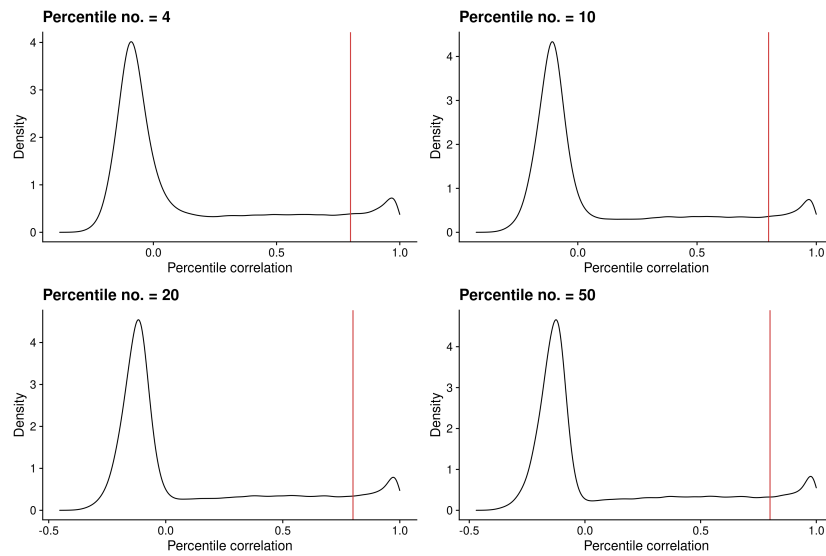
Supplementary Figure 3. A) Heterogeneity and sparsity in the Tasic et al. dataset. x-axis shows number of cells per isoform showing zero expression (yellow) and expression higher than the isoform mean (pink), y-axis shows density. B) Dataset mean-variance relationship. Dots represent transcript isoforms, x-axis represents the isoform mean and y-axis correspond to isoform expression variance (\log_2 counts). C) Percentile correlation density distributions obtained after using different number of percentiles to summarize cell type-level expression.



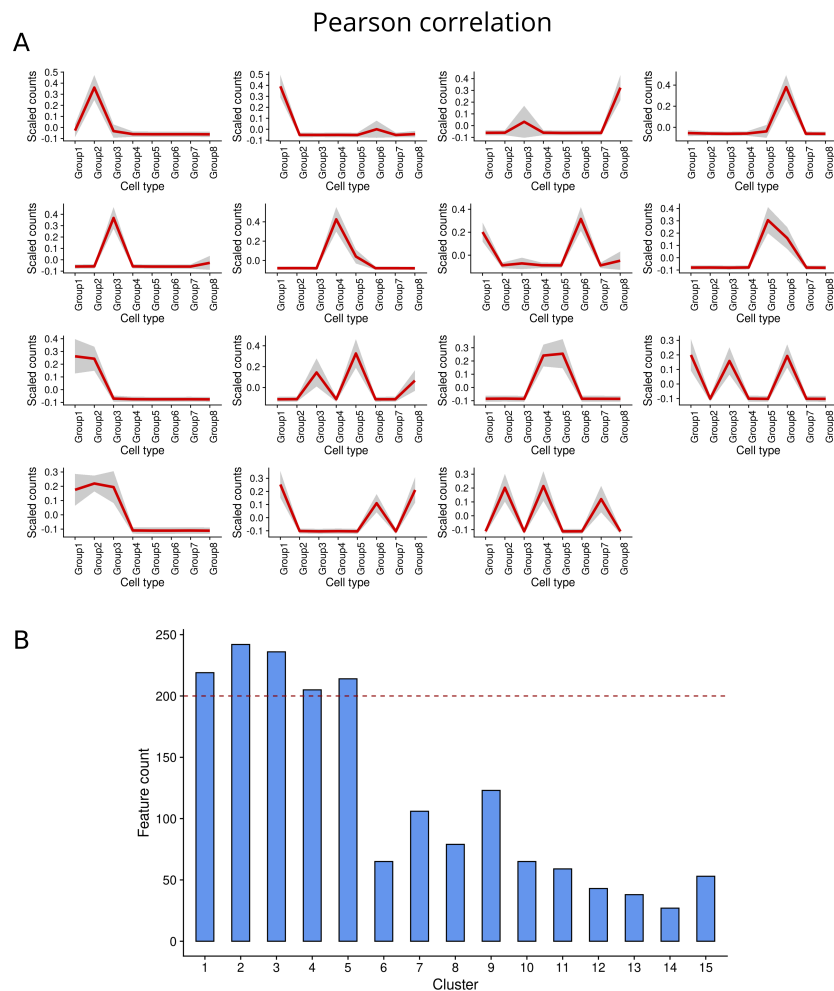
Supplementary Figure 4: cluster refinement examples. A) Example of four redundant clusters, i.e. clusters representing the same expression profiles, merged by profile similarity. B) Example of noisy clusters, i.e. clusters grouping isoforms with highly dissimilar expression patterns across cell types, whose isoforms were re-assigned to other clusters.



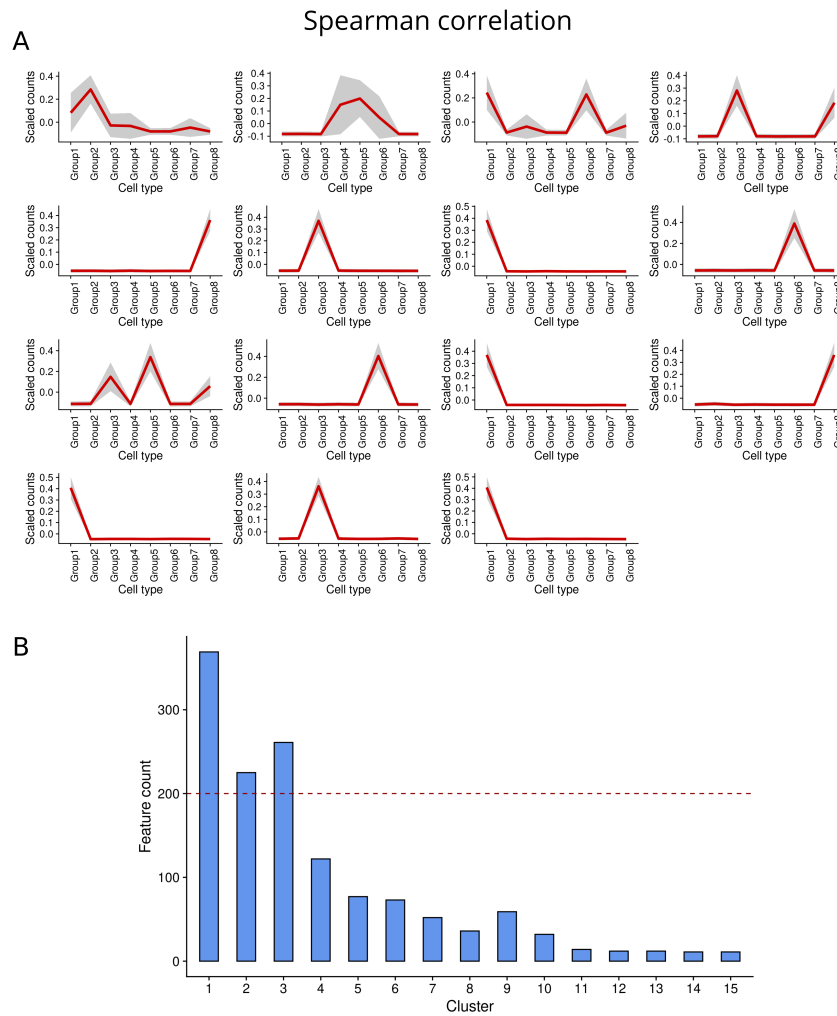
Supplementary Figure 5: simulated data characterization. A) Single-cell co-expression pattern simulation strategy. After simulating single-cell RNA-seq data with SymSim, features are re-ranked within each cell type, using the simulated expression values to generate groups of co-expressed transcripts. Simulated clusters were generated according to 15 pre-defined patterns. B) and C) Expression profile of the 15 simulated clusters obtained, before (B) and after (C) filtering synthetic features showing dissimilarity with the simulated profile. Cell-level mean expression (scaled, see Methods) is computed for all simulated transcripts in the cluster and then aggregated as the global cell type mean, represented by the red line. Grey area corresponds to cell type mean \pm standard deviation. D) Number of transcripts in simulated clusters after filtering to ensure cluster profile consistency. E) Heatmap representing the proportion of high co-expression values (correlation or proportionality > 0.8), computed for all pairs of synthetic transcripts within each simulated cluster. Darker colors indicate successful recapitulation of simulated co-expression relationships by the evaluated co-expression metrics.



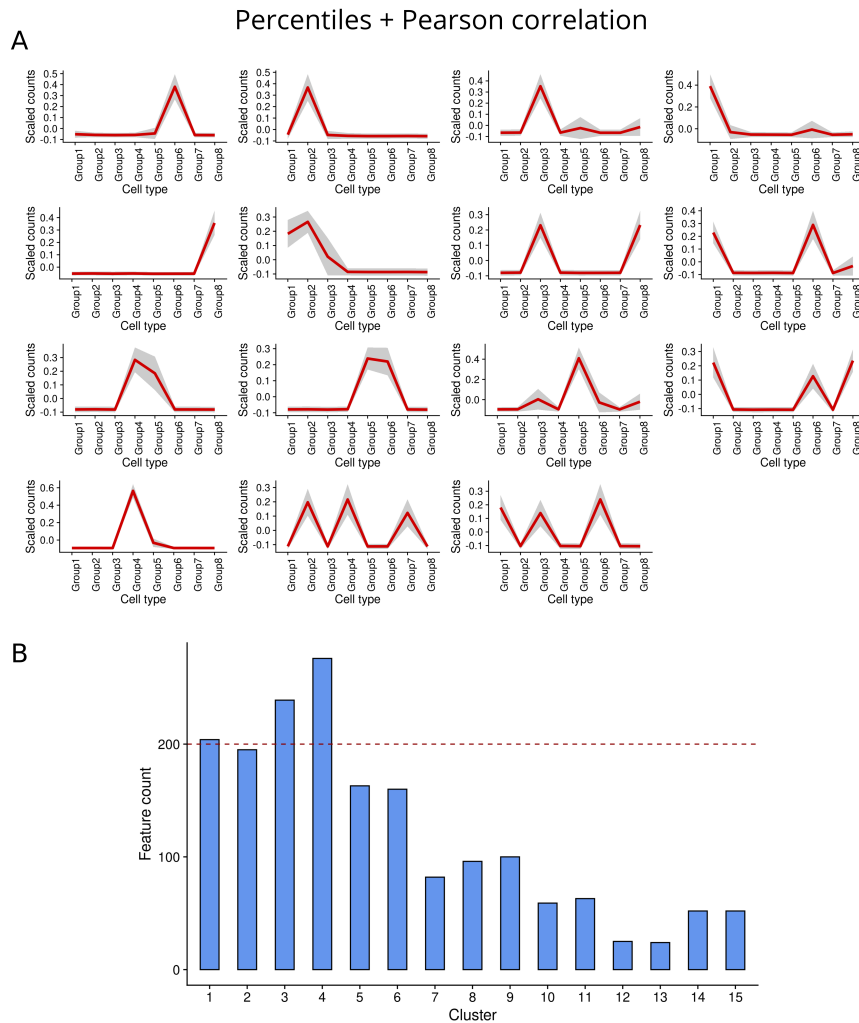
Supplementary Figure 6: effect of percentile number in simulated data. Percentile correlation density distributions obtained after using different number of percentiles to summarize cell type-level expression.



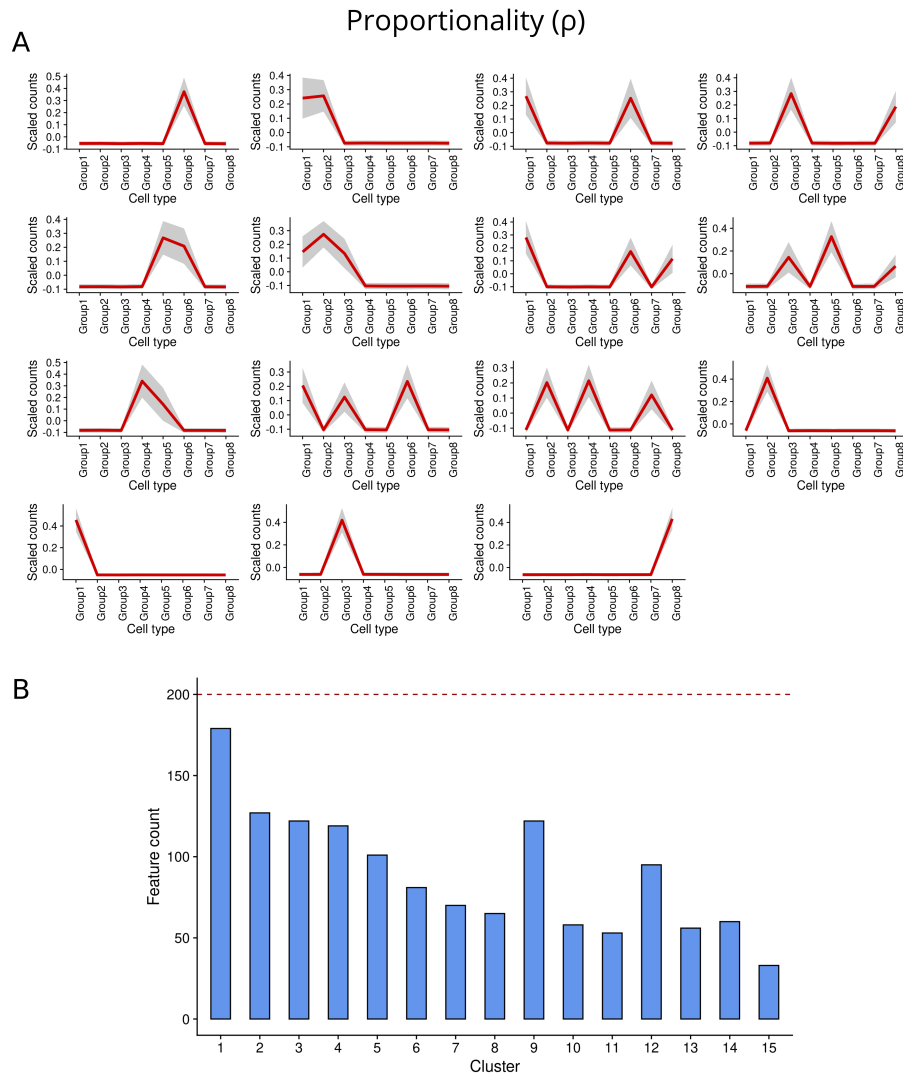
Supplementary Figure 7: clustering of simulated single-cell data using Pearson correlation as distance. A) Mean profiles of generated clusters. Cell-level mean expression (scaled, see Methods) was computed for all simulated transcripts in the cluster and then aggregated as the global cell type mean, represented by the red line. Grey area corresponds to cell type mean \pm standard deviation. B) Number of simulated transcripts per cluster.



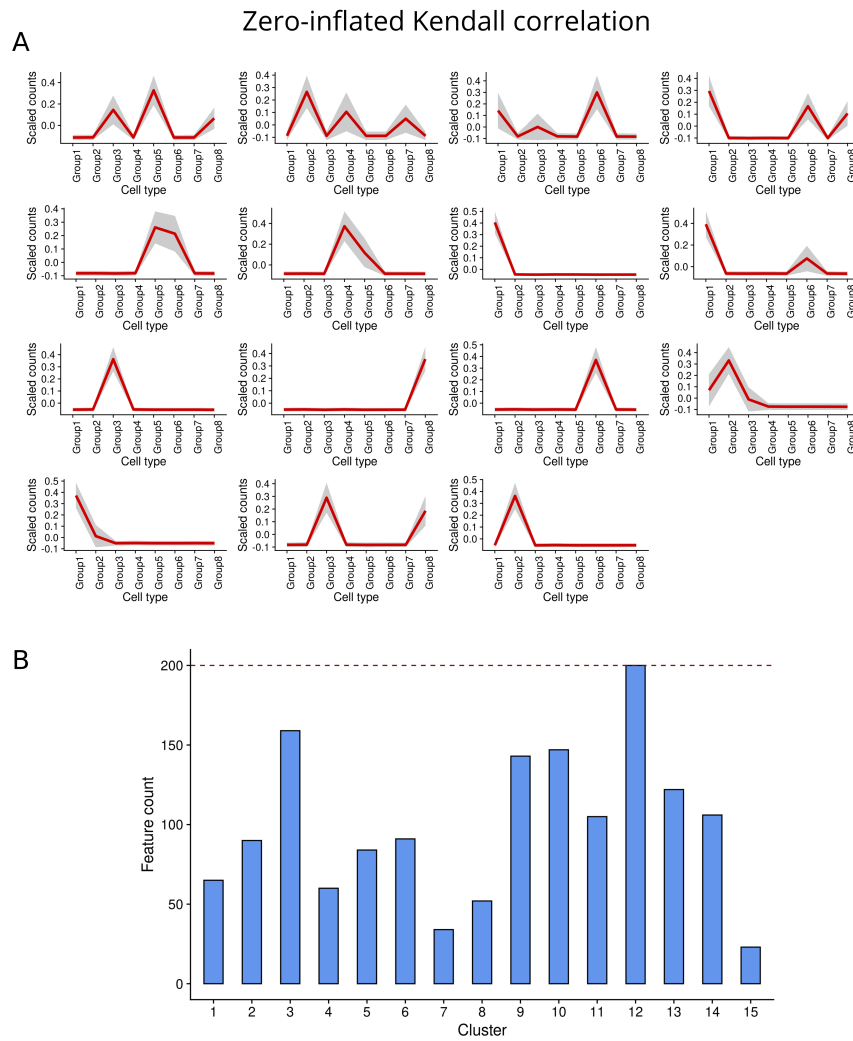
Supplementary Figure 8: clustering of simulated single-cell data using Spearman correlation as distance. A) Mean profiles of generated clusters. Cell-level mean expression (scaled, see Methods) was computed for all simulated transcripts in the cluster and then aggregated as the global cell type mean, represented by the red line. Grey area corresponds to cell type mean \pm standard deviation. B) Number of simulated transcripts per cluster.



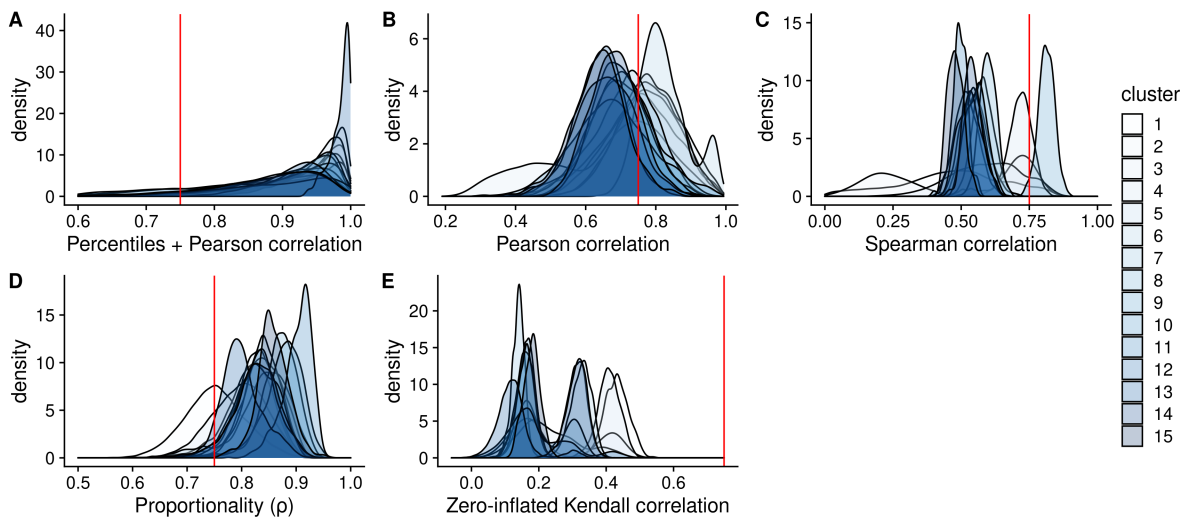
Supplementary Figure 9: clustering of simulated single-cell data using percentiles + Pearson correlation as distance. A) Mean profiles of generated clusters. Cell-level mean expression (scaled, see Methods) was computed for all simulated transcripts in the cluster and then aggregated as the global cell type mean, represented by the red line. Grey area corresponds to cell type mean \pm standard deviation. B) Number of simulated transcripts per cluster.



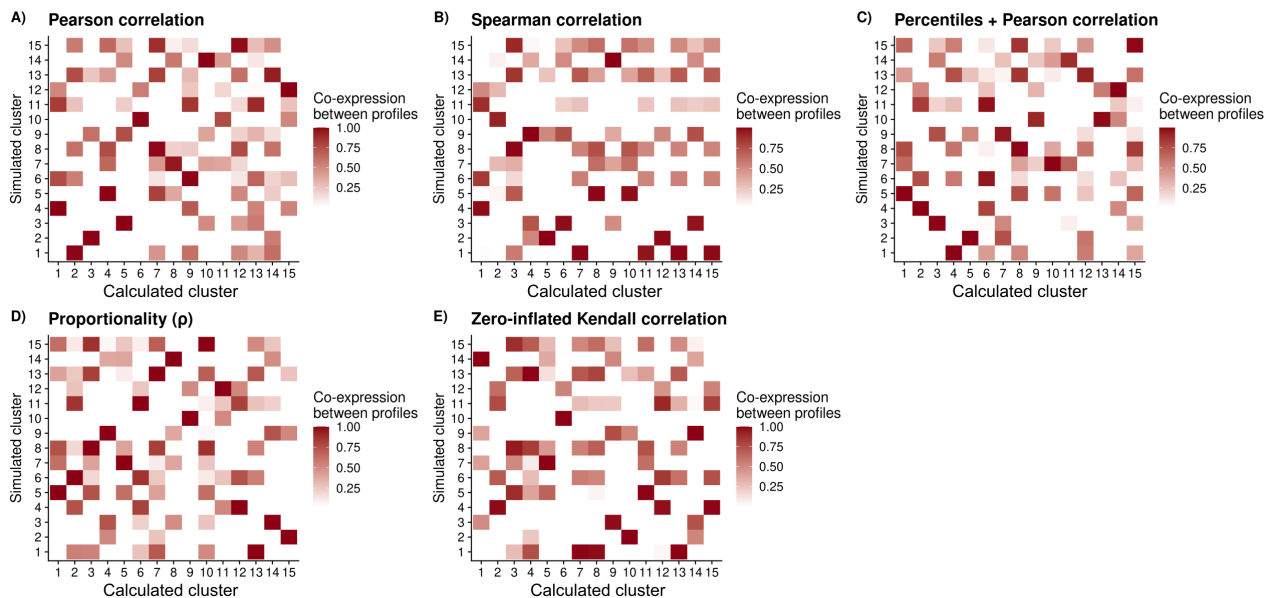
Supplementary Figure 10: clustering of simulated single-cell data using proportionality (ρ) as distance. A) Mean profiles of generated clusters. Cell-level mean expression (scaled, see Methods) was computed for all simulated transcripts in the cluster and then aggregated as the global cell type mean, represented by the red line. Grey area corresponds to cell type mean \pm standard deviation. B) Number of simulated transcripts per cluster.



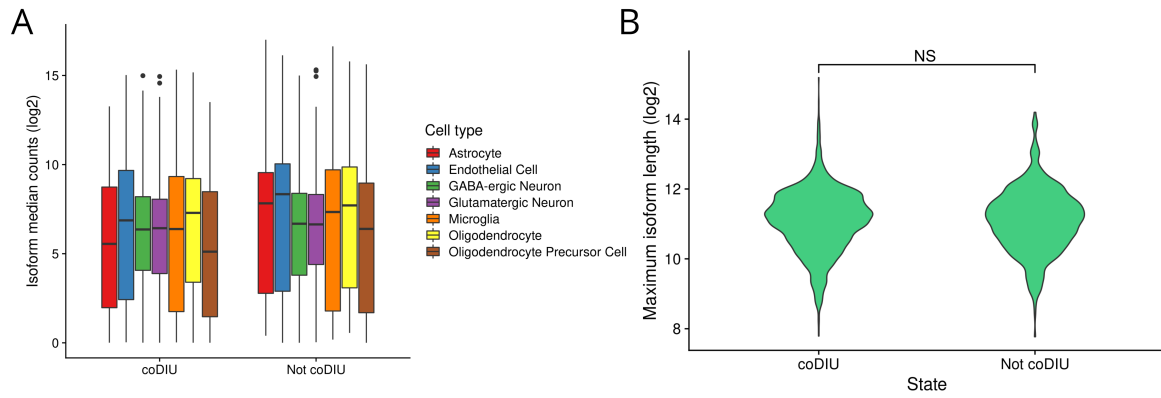
Supplementary Figure 11: clustering of simulated single-cell data using zero-inflated Kendall correlation as distance. A) Mean profiles of generated clusters. Cell-level mean expression (scaled, see Methods) was computed for all simulated transcripts in the cluster and then aggregated as the global cell type mean, represented by the red line. Grey area corresponds to cell type mean \pm standard deviation. B) Number of simulated transcripts per cluster



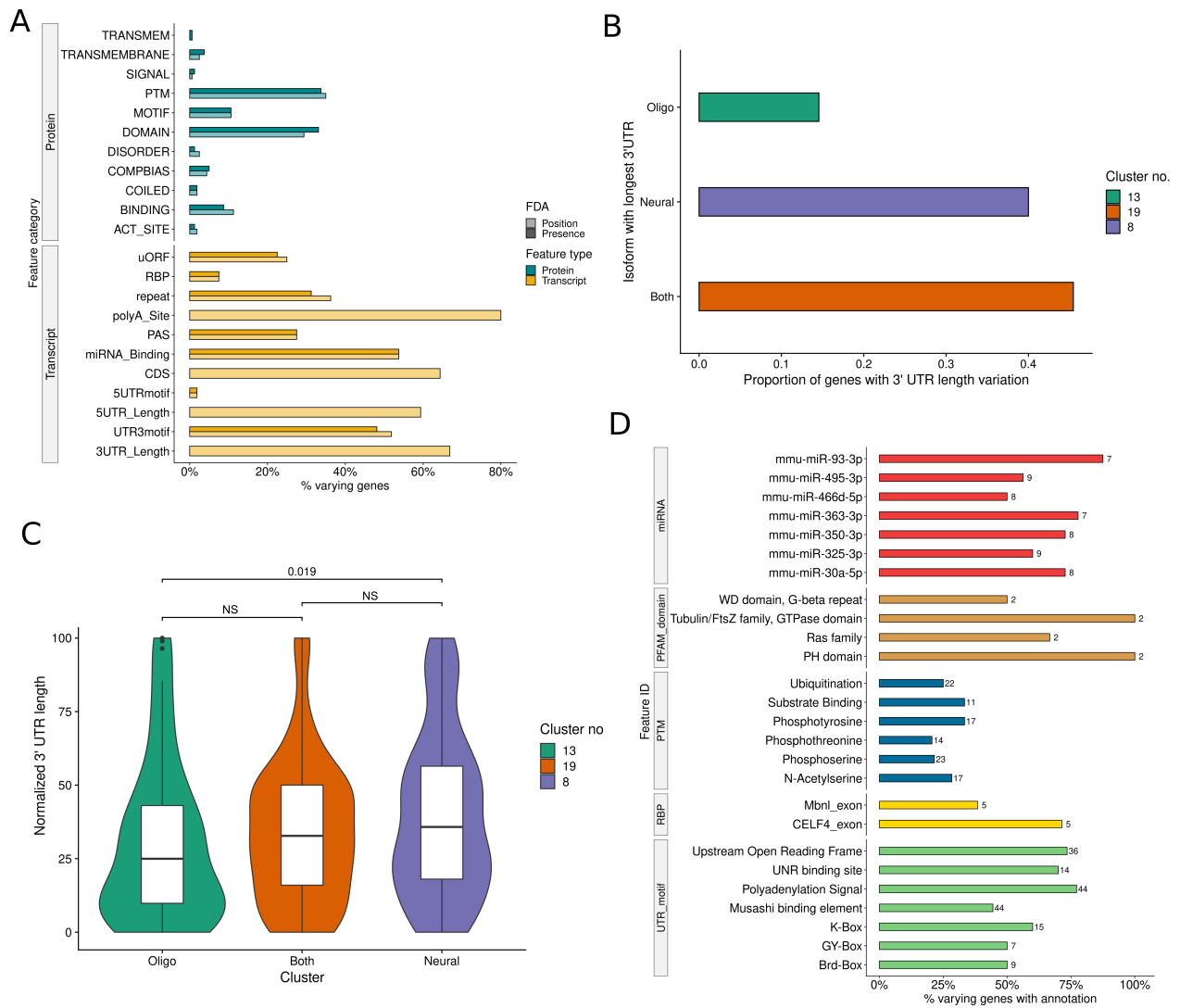
Supplementary Figure 12: density distributions of co-expression values between isoforms in clusters obtained from simulated data. Red line indicates the 0.75 threshold value. Results shown for the clusters generated using each evaluated metric: A) Percentiles + Pearson correlation. B) Pearson correlation. C) Spearman correlation. D) Proportionality (ρ). E) Zero-inflated Kendall correlation.



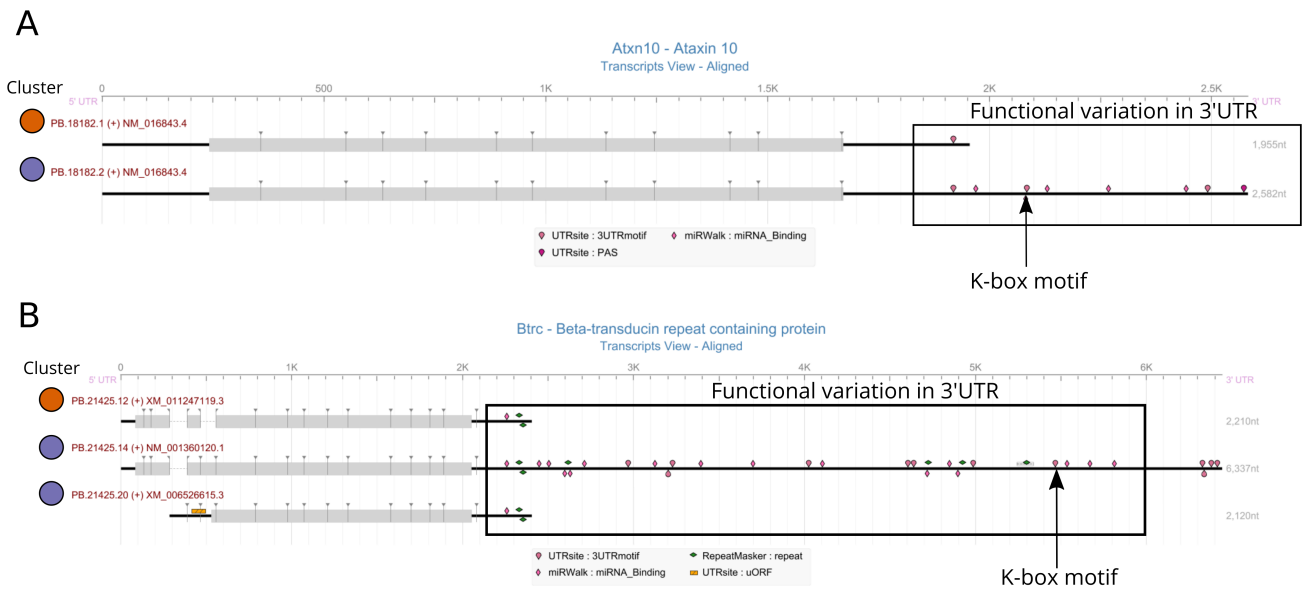
Supplementary Figure 13: Pearson correlation between mean expression profiles of calculated and simulated clusters. Results shown for the clusters generated using each evaluated metric: A) Pearson correlation. B) Spearman correlation. C) Percentiles + Pearson correlation. D) Proportionality (ρ). E) Zero-inflated Kendall correlation. Calculated clusters show no unique match among simulated clusters and were therefore paired with the simulated cluster showing the highest mean expression profile correlation.



Supplementary Figure 14: evaluation of expression and length biases on coDIU genes. A) Boxplot of cell type-level median counts (log₂) for isoforms belonging to coDIU and not coDIU (i.e. solely DIU) genes. Dots represent expression outliers. B) Violin plot of isoform length (log₂) of longest isoform in coDIU and not coDIU (i.e. solely DIU) genes. No significant differences were found between them (NS: not significant, Wilcoxon test).



Supplementary Figure 15: functional analysis of the 160 coDIU genes detected across neural-oligodendrocyte clusters (no. 8, 13 and 19). A) Functional Diversity Analysis (FDA) results. y-x: functional annotation categories. x-axis: percentage of genes including at least one *varying* feature from a given functional category. B) Proportion of coDIU genes with 3'UTR variation across isoforms that have their longest 3'UTR isoform in each of the three analyzed clusters. C) Violin + boxplot of normalized 3'UTR lengths for each of the three neural-oligodendrocyte clusters. Normalized lengths are computed by dividing each individual isoform's 3'UTR length by the sum of 3'UTR lengths of all the gene's isoforms. Significance levels for the comparison of the three groups are indicated above the corresponding braces (*p*-value, Wilcoxon test). D) ID-level FDA results for features in highly-varying functional categories. Percentage of genes containing the feature that show functional variation across isoforms (x-axis) is shown for the most frequently annotated features in each category. Total coDIU genes with feature indicated by the bar label.



Supplementary Figure 16: tappAS view of transcript functional annotation for A) *Atxn10* and B) *Btrc* isoforms. Cluster assignments for each isoform are indicated by dot color.