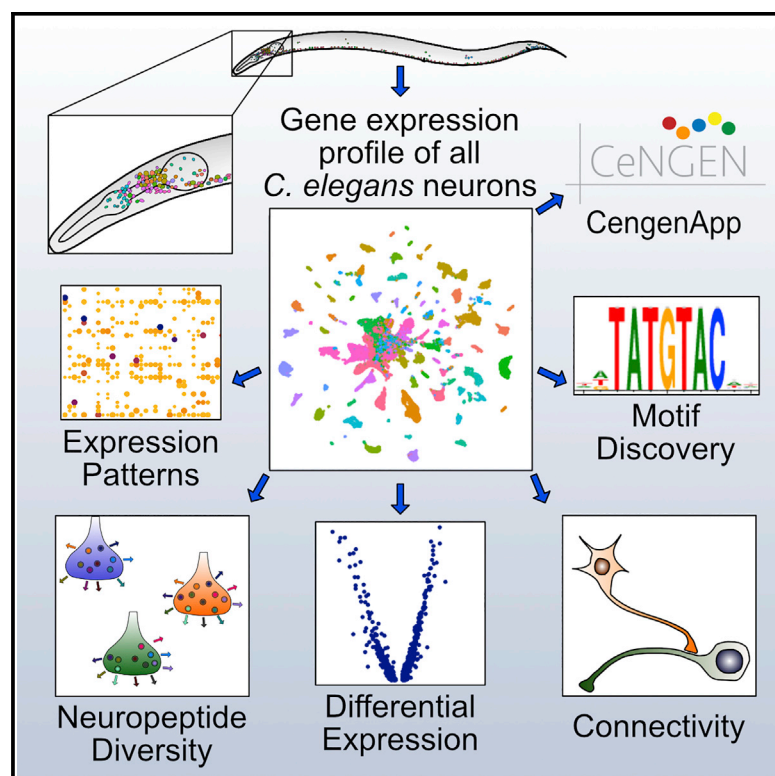


# Molecular topography of an entire nervous system

## Graphical abstract



## Authors

Seth R. Taylor, Gabriel Santpere, Alexis Weinreb, ..., Marc Hammarlund, Oliver Hobert, David M. Miller III

## Correspondence

marc.hammarlund@yale.edu (M.H.),  
or38@columbia.edu (O.H.),  
david.miller@vanderbilt.edu (D.M.M.)

## In brief

A gene expression map captures all 302 neurons in mature *C. elegans* deciphering the molecular basis for cell heterogeneity, connectivity, and function.

## Highlights

- Gene expression profiles of all 118 neuron classes in the *C. elegans* hermaphrodite
- Each neuron type expresses a distinct code of neuropeptide genes and receptors
- Expression profiles enable discovery of cell-type-specific *cis*-regulatory sequences
- Cell adhesion molecules correlate with neuron-specific connectivity



## Resource

# Molecular topography of an entire nervous system

Seth R. Taylor,<sup>1</sup> Gabriel Santpere,<sup>2,3,10</sup> Alexis Weinreb,<sup>2,4,10</sup> Alec Barrett,<sup>2,4,10</sup> Molly B. Reilly,<sup>5,6,10</sup> Chuan Xu,<sup>2,10</sup> Erdem Varol,<sup>7,10</sup> Panos Oikonomou,<sup>5,8,10</sup> Lori Glenwinkel,<sup>5,6</sup> Rebecca McWhirter,<sup>1</sup> Abigail Poff,<sup>1</sup> Manasa Basavaraju,<sup>2,4</sup> Ibnul Rafi,<sup>5,6</sup> Eviatar Yemini,<sup>5,6</sup> Steven J. Cook,<sup>5,6</sup> Alexander Abrams,<sup>2,4</sup> Berta Vidal,<sup>5,6</sup> Cyril Cros,<sup>5,6</sup> Saeed Tavazoie,<sup>5,8</sup> Nenad Sestan,<sup>2</sup> Marc Hammarlund,<sup>2,4,\*</sup> Oliver Hobert,<sup>5,6,\*</sup> and David M. Miller III<sup>1,9,11,\*</sup>

<sup>1</sup>Department of Cell and Developmental Biology, Vanderbilt University School of Medicine, Nashville, TN, USA

<sup>2</sup>Department of Neuroscience, Yale University School of Medicine, New Haven, CT, USA

<sup>3</sup>Neurogenomics Group, Research Programme on Biomedical Informatics (GRIB), Hospital del Mar Medical Research Institute (IMIM), DCEXS, Universitat Pompeu Fabra, 08003 Barcelona, Catalonia, Spain

<sup>4</sup>Department of Genetics, Yale University School of Medicine, New Haven, CT, USA

<sup>5</sup>Department of Biological Sciences, Columbia University, New York, NY, USA

<sup>6</sup>Howard Hughes Medical Institute, Columbia University, New York, NY, USA

<sup>7</sup>Department of Statistics, Columbia University, New York, NY, USA

<sup>8</sup>Department of Systems Biology, Columbia University Medical Center, New York, NY, USA

<sup>9</sup>Program in Neuroscience, Vanderbilt University School of Medicine, Nashville, TN, USA

<sup>10</sup>These authors contributed equally

<sup>11</sup>Lead contact

\*Correspondence: [marc.hammarlund@yale.edu](mailto:marc.hammarlund@yale.edu) (M.H.), [or38@columbia.edu](mailto:or38@columbia.edu) (O.H.), [david.miller@vanderbilt.edu](mailto:david.miller@vanderbilt.edu) (D.M.M.)

<https://doi.org/10.1016/j.cell.2021.06.023>

## SUMMARY

We have produced gene expression profiles of all 302 neurons of the *C. elegans* nervous system that match the single-cell resolution of its anatomy and wiring diagram. Our results suggest that individual neuron classes can be solely identified by combinatorial expression of specific gene families. For example, each neuron class expresses distinct codes of ~23 neuropeptide genes and ~36 neuropeptide receptors, delineating a complex and expansive “wireless” signaling network. To demonstrate the utility of this comprehensive gene expression catalog, we used computational approaches to (1) identify *cis*-regulatory elements for neuron-specific gene expression and (2) reveal adhesion proteins with potential roles in process placement and synaptic specificity. Our expression data are available at <https://cengen.org> and can be interrogated at the web application CengenApp. We expect that this neuron-specific directory of gene expression will spur investigations of underlying mechanisms that define anatomy, connectivity, and function throughout the *C. elegans* nervous system.

## INTRODUCTION

Neurons share many common functions, yet there are a remarkable variety of different neuronal types, each with distinct features and functions. Because genetic programs likely specify these differences, a comprehensive molecular model of the brain requires a gene expression map at single-cell resolution. Although profiling methods have cataloged diverse neuron types in a variety of organisms (Adorjan et al., 2019; Poulin et al., 2016; Tasic et al., 2016; Zeisel et al., 2015; Zhu et al., 2018), incomplete knowledge of the anatomy and wiring of complex nervous systems has hampered the effort to link neuron-specific functional and anatomical properties with individual molecular signatures.

To investigate the relationship between gene expression and neuroanatomy, we produced single-cell RNA sequencing (scRNA-seq) profiles for all neuron types in an entire nervous system, that of the *C. elegans* hermaphrodite. The complete anatomy and wiring diagram of the *C. elegans* nervous system were defined by serial section electron microscopy (Albertson

and Thomson, 1976; Brittin et al., 2021; Cook et al., 2019; White et al., 1986; Witvliet et al., 2020). This approach identified 118 anatomically distinct classes among the 302 neurons in the mature hermaphrodite nervous system. We established the *C. elegans* Neuronal Gene Expression Map & Network (CeN-GEN) consortium (Hammarlund et al., 2018) to generate transcriptional profiles of each neuron class, thereby bridging the gap between *C. elegans* neuroanatomy and the genetic blueprint that defines it. We used fluorescence activated cell sorting (FACS) to isolate neurons from L4 stage larvae for scRNA-seq. By the L4 stage, the entire nervous system has been generated and most neurons have terminally differentiated. Our approach generated profiles of 70,296 neurons, including all 118 canonical neuron classes and thus offers a comprehensive catalog of gene expression for an entire nervous system.

We found that every neuron class is defined by distinct combinations of neuropeptide-encoding genes and neuropeptide receptors, suggesting different roles for each type of neuron in sending and receiving signals. We identified an expansive



catalog of DNA and RNA sequence motifs that are correlated with cohorts of co-regulated genes. We used computational approaches to identify cell adhesion molecules associated with neuron-specific synapses and bundling. Together, our results provide a comprehensive link between neuron-specific gene expression and the structure and function of an entire nervous system. We expect that these datasets and the tools that we have developed for interrogating them will power future investigations into the genetic basis of neuronal connectivity and function.

## RESULTS AND DISCUSSION

### scRNA-seq identifies all known neuron classes in the mature *C. elegans* nervous system

To profile the entire *C. elegans* nervous system (Figure 1A), we isolated neurons at the L4 larval stage, when all neuron types have been generated (Sulston and Horvitz, 1977) and terminally differentiated to generate a functional nervous system. Initially, we used FACS to isolate neurons from a pan-neural marker strain and found that many neuron classes were either underrepresented or absent (Figures S1A–S1C). To overcome this limitation, we isolated cells from a series of fluorescent marker strains that labeled distinct subsets of neurons (Figure 1C; Table S1). We generated 100,955 single cell transcriptomes with a median of 928 unique molecular identifiers (UMIs) and 328 genes/cell. Application of the uniform manifold approximation and projection (UMAP) dimensional reduction algorithm effectively segregated most of these cells into distinct groups (Figure S2A).

We separated non-neuronal cells (27,427 cells, 27.2%) (Figures S2B–S2F) and neurons (70,296 cells, 69.6%) (Figures 1A and 1B) into different sub-UMAPs for further annotation. Neurons had a median of 1,033 UMIs and 363 genes/cell. Most neuronal UMAP clusters could be assigned to individual neuron classes based on known marker genes (Hobert et al., 2016; Reilly et al., 2020; Figures S3A–S3C). For clusters that could not be so readily identified, we generated GFP transcriptional reporters for genes enriched in the target clusters for direct examination *in vivo* (Figures 1D, S3D, and S3E). For example, *C39H7.2* was exclusively detected in a small cluster that expressed no known distinct markers. We used the multi-colored NeuroPAL marker strain (Yemini et al., 2021) to determine that a *C39H7.2::NLS-GFP* transcriptional reporter was exclusively expressed in the tail interneuron LUA (Figure 1D).

Ninety of the 118 neuronal types were detected in distinct clusters in the pan-neuronal UMAP (Figure 1B). The remaining clusters contained multiple, closely related neuron classes (e.g., oxygen-sensing neurons, ventral cord motor neurons). Individual UMAP projections of these clusters facilitated the annotation of 38 additional neuron types (Figures 1E, 1F, and S3F–S3J), including subtypes within 10 classes (see below). Only two neuron classes were inseparable, the DD and VD ventral cord GABAergic motor neurons, despite known differences in gene expression (Mekman and Sengupta, 2005; Petersen et al., 2011; Shan et al., 2005). Overall, we annotated 95.9% of the cells in the entire dataset and identified distinct clusters encompassing all of the 118 anatomically defined neuron classes in the mature hermaphrodite nervous system (White et al., 1986).

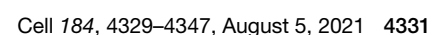
### scRNA-seq reveals transcriptionally distinct neuronal subtypes

Reporter-based gene expression and connectivity data suggest that some of the 118 anatomically defined neuron classes may be comprised of separate subclasses (Hobert et al., 2016; White et al., 1986). Our results confirmed this prediction by revealing 128 transcriptionally distinct neuron types, including subtypes within 10 of the 118 canonical neuron classes. Consistent with earlier findings (Cao et al., 2017; Johnston et al., 2005; Lesch et al., 2009; Packer et al., 2019; Pierce-Shimomura et al., 2001; Troemel et al., 1999; Vidal et al., 2018; Yu et al., 1997), we detected individual clusters for the bilaterally asymmetric sensory neuron pairs ASE (ASER and ASEL) and AWC (AWC<sup>ON</sup> and AWC<sup>OFF</sup>) (Figures 2A and S4A). Differential gene expression analysis revealed expanded lists of subtype-specific transcripts for the ASE and AWC subclasses (Figures 2B and S4B), including asymmetric expression of receptor-type guanylyl cyclases (rGCs) (Ortiz et al., 2006) and neuropeptides (Figures 2A, 2B, and S4A). Other than the AWC and ASE neuron pairs, we detected no other cases of molecularly separable left/right homologous cells within a neuron class.

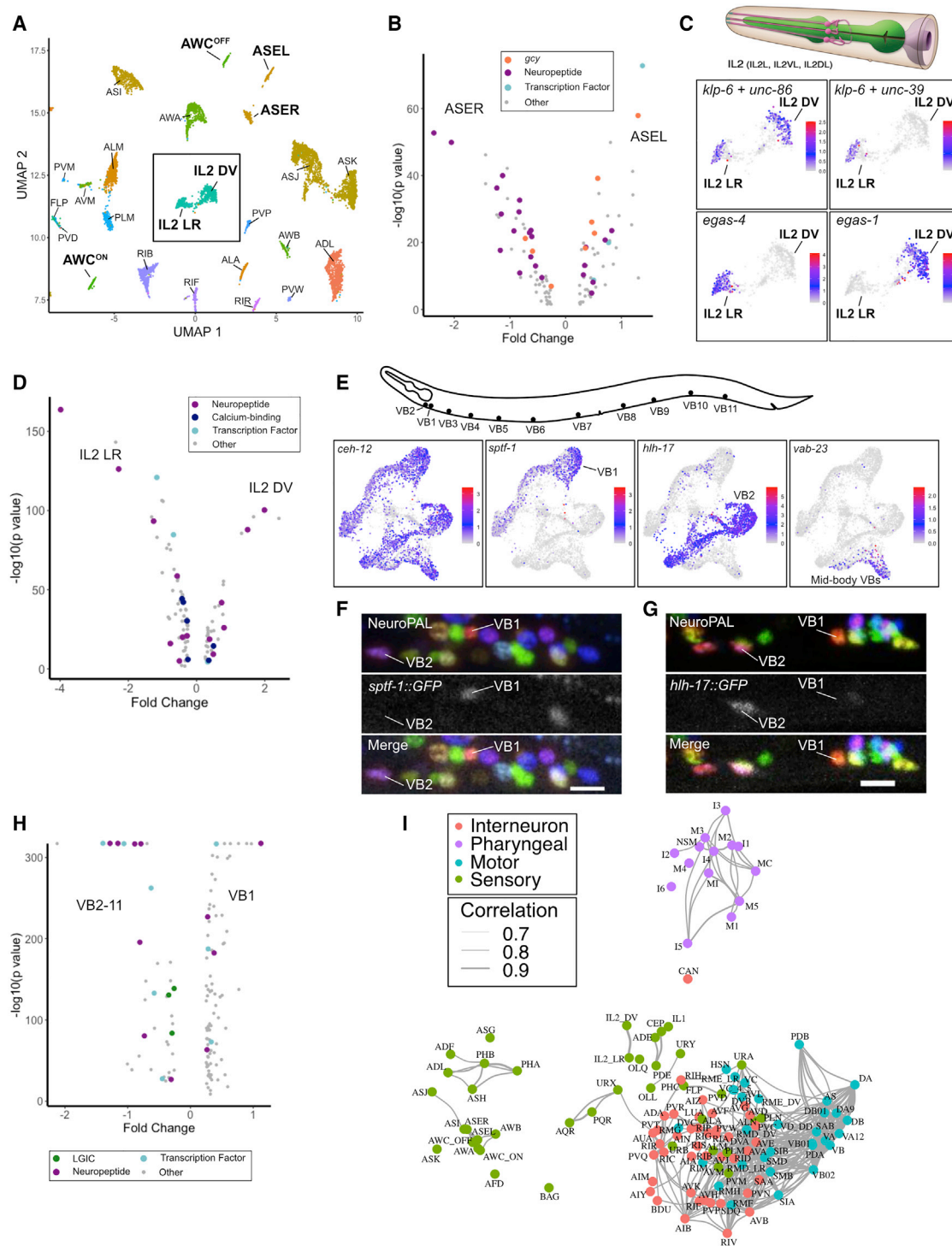
The remaining eight neuron classes with transcriptionally distinct subtypes are either arranged in radially symmetric groups of 4 or 6 neurons or are distributed along the anterior/posterior axis in the motor circuit. We detected distinct subclusters for two neuron classes with 6-fold symmetry at the nerve ring, the inner labial IL2 neurons (Figures 2A and 2C) and the RMD neurons (Figures 1E and S4A). In both cases, the left/right pair of neurons (e.g., IL2L/R) segregates from the dorsal/ventral pairs (IL2DL/R and IL2VL/R). Differentially expressed genes between the IL2 clusters encode neuropeptides, ion channels, calcium binding proteins, and transcription factors and point to potentially distinct functions for the subtypes (Figures 2C and 2D). For the GABAergic RME head motor neurons, we detected distinct dorsal/ventral (RMED/V) and left/right clusters (RMEL/R) (Figures 1F and S4A). We also identified multiple clusters for the DA, DB, VA, VB, and VC ventral nerve cord motor neuron classes. In each case, one subtype corresponded to one or two individual members of these classes. For example, VC4 and VC5, which flank the vulva, clustered independently from the other four VC neurons (Figures 1F and S4A). For A-class motor neurons (DA and VA), we detected distinct clusters corresponding to the most posterior neurons located in the pre-anal ganglion, DA9 and VA12 (Figures 1E and S4A).

Both B-class motor neuron classes (DB and VB) contained multiple independent clusters (Figures 2E and S4A). In this case, the most anterior B-class motor neurons (DB1, VB1, and VB2) segregated into separate clusters. The homeodomain transcription factor CEH-12 is selectively expressed in VBs (Von Steina et al., 2007) and marks the VB clusters (Figure 2E). We identified VB1 based on expression of a GFP reporter gene for the subcluster-specific marker *sptf-1* (Figure 2E–F). The VB2 subcluster was similarly identified by the selective expression of *hlh-17::GFP* in VB2 among VBs *in vivo* (Figure 2E–G). Interestingly, all of the molecularly distinct subclasses we detected also have known differences in synaptic connectivity (Hobert et al., 2016; White et al., 1986).

We did not detect subtypes for additional classes with 3-, 4-, or 6-fold symmetry. This may be due to the low number of cells







**Figure 2. Identification of neuron subtypes**

(A) UMAP of neurons with molecularly distinct subtypes (bold labels) from neuronal UMAP (Figure 1B). Inset denotes IL2 DV and IL2 LR clusters.

(B) Volcano plot of differentially expressed genes (false discovery rate [FDR] <0.05) for ASER versus ASEL. Guanylyl cyclases (*gcy*), neuropeptides, and transcription factors are marked.

(C) Top: 3 pairs of IL2 sensory neurons (IL2L/R, IL2VL/R, and IL2DL/R) from WormAtlas. Bottom: UMAP inset from (A) showing normalized expression of marker genes for all IL2 neurons (*klp-6*, *unc-86*), IL2 LR (*unc-39*, *egas-4*), and IL2 DV (*egas-1*).

(D) Volcano plot of differentially expressed genes (FDR <0.05) between IL2 subtypes.

(legend continued on next page)

(<100 for OLQ, SAA, URY, and IL1, see [Table S1](#)) assigned to some of these classes. Alternatively, molecular differences among subsets of these neuron types ([Hobert et al., 2016](#)) may be limited to a small number of genes that would be insufficient to drive separation in our analyses.

Using 7,390 highly variable genes (see [STAR Methods](#)), we generated a network describing the relative molecular relationship of the 128 identified neuron classes and subclasses ([Figure 2I](#)). This approach separated sensory and motor neurons as well as a distinct cluster of pharyngeal neurons. Interestingly, pre-motor interneurons cluster with motor neurons. Amphid/phasmid sensory neurons clearly separated from non-amphid/phasmid sensory neuron types. Within amphid/phasmid neurons, some neurons cluster according to sensory modalities. Notably, the chemorepulsive neurons ADL, ASH, and PHA/PHB form their own subcluster. The CO<sub>2</sub>-sensitive BAG neuron and the CAN neuron show the least similarity to other neuron types. Thus, a systematic comparison of neuron-specific profiles confirms that neurons with shared anatomical and functional characteristics are defined by similar patterns of gene expression.

### Defining gene expression across neuron types

A key consideration for scRNA-seq data is accurately determining whether a detected signal (UMI) for a given gene is actual expression in a cell type (rather than noise). We addressed this question quantitatively by thresholding aggregated data for each cell type using a ground-truth dataset of high-confidence gene expression results across the entire nervous system (mostly fosmid-based reporters and/or reporter-tagged endogenous genes; see [STAR Methods](#) and [Figure S5](#)). We selected 4 threshold levels (designated as 1–4) offering different compromises between the risk of false-positives and false-negatives. We used threshold 2 for subsequent analyses. With this threshold, we estimate a true positive detection rate of 0.81 and a false discovery rate of 0.14 (see [STAR Methods](#)). The number of genes detected per neuron type (median, 5,842; range = 1,371 [ALN] to 7,542 [ASJ]) was positively correlated with the number of cells sequenced per neuron type (median, 352; range = 12 [M4] to 3,189 [AIJ]; [Figure S5I](#), Spearman rank correlation = 0.783,  $p < 2.2 \times 10^{-16}$ ) and with the true positive rate ([Figure S5J](#), Spearman rank correlation = 0.6776,  $p < 2.2 \times 10^{-16}$ ). Neurons with fewer cells and fewer detected genes were concentrated in the anterior and pre-anal ganglia ([Figure S5H](#)), possibly reflecting bias in the dissociation procedure. Nine neuron classes with the fewest detected genes and lowest true positive rates compared to ground truth are labeled in [Figure S5J](#). These cell types are likely to have the highest rates of false negatives, as we estimate the true mean number of genes expressed per neuron type to be ~6,550 (see [STAR Methods](#)).

We examined the distribution of genes encoding ribosomal proteins to test whether our thresholding approach would preserve a predicted ubiquitous pattern of gene expression. Our results show that 65 of the 78 ribosomal genes (83%) are detected in  $\geq 98\%$  of neuron types, with 53 (68%) expressed in all but one cell type (ALN) ([Figure 3A](#)). Overall, these results indicate that our thresholding approach accurately identifies expressed genes for most cell types in the *C. elegans* nervous system.

### Neuron-specific codes of neuropeptide signaling genes

We used the thresholded dataset (threshold 2) to probe expression of selected gene families known to be involved in various aspects of neuron function and development ([Data S1](#)) and provide highlights of this analysis here in the main text. Neuropeptide-encoding genes (31 FMRFamide-like peptides [*flp*], 33 insulin-related peptides [*ins*], and 77 neuropeptide-like proteins [*nlp*] genes, total of 141 genes) were detected in every neuron class (a minimum of 6, maximum of 62 per neuron) ([Figure 3](#)). Consistently, neuropeptide processing genes were broadly expressed throughout the nervous system ([Figure 3A](#)). Strikingly, each neuron class expressed a distinct combination of neuropeptides, averaging 23 genes. Sensory neurons and interneurons expressed more neuropeptide genes than motor neurons ([Figure 3E](#)). Further, neuropeptide encoding genes are among the most highly expressed transcripts in our dataset, similar to reports from *Hydra*, *Drosophila*, and mouse neurons ([Siebert et al., 2019](#); [Allen et al., 2020](#); [Smith et al., 2019](#)). Moreover, the subset of 25 *nlp* genes with homologs in other species ([Husson et al., 2009](#); [Kozioł et al., 2016](#); [Mirabeau and Joly, 2013](#)), along with the *flp* family genes, were detected at higher levels than *ins* and non-conserved *nlp* genes ([Figure 3B](#)).

Whereas several neuropeptide-encoding genes (*flp-9*, *flp-5*, and *nlp-21*) were widely expressed, we also detected neuropeptides with expression restricted to just one or two neuron types, including exclusive expression of *flp-1* in AVK, *flp-23* in HSN, *nlp-56* in RMG, *nlp-2* and *nlp-23* in AWA and *ins-13* in RMED/V ([Figure 3C](#)). We validated the restricted expression of *nlp-56* in the RMG cluster and *flp-1* in AVK with CRISPR/Cas9-engineered reporter alleles ([Figure 3D](#); see also [Figure S6](#)).

Of the more than 140 neuropeptide receptors, most show highly restricted expression, with a few notable exceptions ([Figure 3A](#)). The predicted neuropeptide receptors *pdf-1*, *npr-23*, and *F59D12.1* were expressed in over 100 neuron types. *daf-2*, the only insulin/insulin growth factor (IGF) receptor-like tyrosine kinase in *C. elegans*, was detected in 103 of 128 neuron types. Most other neuropeptide receptor genes were expressed in a restricted subset of neurons; half were expressed in 29 or fewer cell types ([Figure 3A](#)). Each individual neuron type expressed a distinct set of neuropeptide receptors, averaging 36 genes. Sensory neurons and interneurons expressed more neuropeptide

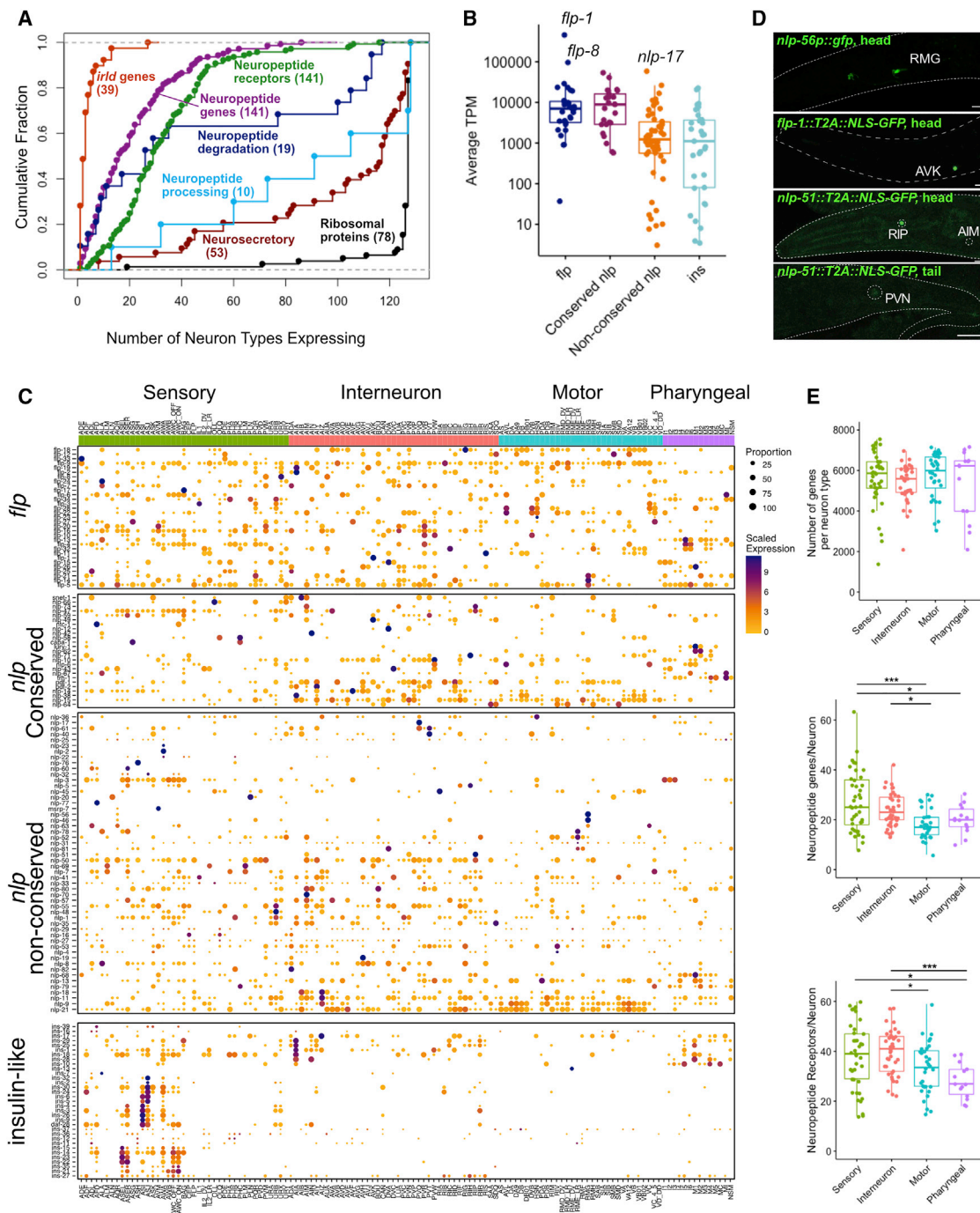
(E) Top: VB motor neuron soma in the ventral nerve cord. Bottom: sub-UMAPs of VB neurons highlighting VB marker (*ceh-12*) and genes (*sptf-1*, *hlh-17*, *vab-23*) expressed in specific VB sub-clusters.

(F and G) Confocal images in NeuroPAL show *sptf-1::GFP* expression in VB1 but not VB2 (F) and selective expression of *hlh-17::GFP* in VB2 but not VB1 (G). Scale bars, 10  $\mu$ m.

(H) Volcano plot of differentially expressed genes (LGICs -ligand-gated ion channels) (FDR <0.05) for VB1 versus all other VB neurons.

(I) *C. elegans* neuron types in a force-directed network by transcriptomic similarities. Colors denote distinct neuron modalities and widths of edges (Pearson correlation coefficients >0.7) show strengths of transcriptome similarity between each pair of neuron types.

See also [Figures S1](#) and [S4](#) and [Table S1](#).



**Figure 3. Expression of neuropeptide signaling genes**

(A) Cumulative distribution plot of neuron types expressing different classes of neuropeptide signaling genes. Each dot is a gene, genes expressed in the same number of neuron types overlap. Numbers in parentheses denote the sum of genes in each category.

(B) Average expression (TPM) for neuropeptide subfamilies across neuron types. *flp-1*, *flp-8*, *nlp-17* are highly expressed. Boxplot spans 25<sup>th</sup> percentile, median and 75<sup>th</sup> percentile.

(C) Heatmap (rows) for *flp* (FMRFamide-related peptide), *nlp* (neuropeptide-like protein), and *ins* (insulin-like peptide) subfamilies across 128 neuron types (columns) grouped by functional/anatomical modalities (sensory, interneuron, motor, pharyngeal). Conserved *nlp* genes are shown separately. Rows are clustered within each family. Circle diameter denotes the proportion of neurons in each cluster that expresses a given gene.

(D) GFP reporters confirm selective expression of *nlp-56* (promoter fusion) in RMG, *flp-1* (CRISPR reporter) in AVK, and *nlp-51* (CRISPR reporter) in RIP, with weaker expression in PVN and AIM. Scale bars, 10  $\mu$ m.

(legend continued on next page)



receptor genes than pharyngeal neurons (Figure 3E). With ongoing efforts to match neuropeptide G protein-coupled receptors (GPCRs) to their cognate ligands (<https://worm.peptide-gpcr.org/project/>), these expression data for all neuropeptide genes and receptors provide a basis for establishing a nervous-system wide map of modulatory neuropeptide signaling.

Signaling complexity across the nervous system is also determined by diverse ionotropic neurotransmitter receptor expression. Each neuron expresses on average 20 ionotropic neurotransmitter receptors, and each individual neuron type expresses a distinct combination of these genes (Data S1). The expression pattern of ionotropic neurotransmitter receptors also suggests extensive non-synaptic volume transmission (Gendrel et al., 2016), further illustrating the complexity of information flow in the *C. elegans* nervous system. The tunability of individual *C. elegans* neurons is illustrated by the wide-spread and complex expression of potassium channels (Data S1). For example, each individual neuron expresses 1 to 18 distinct two-pore TWK-type ion channels.

#### Differential expression of gene regulatory factors

We interrogated gene families involved in gene regulation, including all predicted transcription factors (TFs) (wTF 3.0) (Fuxman Bass et al., 2016) and RNA-binding proteins (Tamburino et al., 2013; Figures 4A–4C; Data S1). 705 of 941 (75%) of predicted transcription factors and 497 of 587 (86%) of predicted RNA-binding proteins were detected in at least one neuron type. Overall, transcription factors were more restricted in their expression than RNA-binding proteins (Figure 4C).

We analyzed expression of all TF classes that contain more than 15 members (homeodomain, nuclear hormone receptor [*nhx*], helix-loop-helix [*bHLH*], C2H2 zinc finger, bZIP, AT hook, and T-box genes) and found distinct themes for individual gene families. At one extreme are T-box genes, only two of which are expressed in postembryonic neurons (Data S1). In contrast, AT hook and bZIP genes are expressed broadly throughout the nervous system. Individual bHLH and C2H2 TF genes show a combination of broad and selective expression in the nervous system (Figure 4C). Each neuron expressed multiple different *nhx* TFs, but sensory and pharyngeal neurons expressed many more *nhx* TFs than either motor neurons or interneurons (Figures 4A–4D). Each amphid and phasmid sensory neuron expressed more than 90 *nhx* TFs. Notably, ASJ expressed 144 *nhx* TFs, 75% of the 191 *nhx* TFs detected in the entire neuronal dataset (Figures 4A and 4B). Abundant expression of a broad array of *nhx* genes in sensory neurons is suggestive of specific roles in mediating transcriptional responses to sensory stimuli.

Homeobox gene expression profiles are distinct from that of other TF families. In agreement with a recent report (Reilly et al., 2020), the majority of homeodomain TFs are sparsely expressed in the nervous system. Most individual homeodomain TFs are selectively expressed in subsets of neuron classes (Figures 4A and 4B). In addition, each neuron class expressed a unique combination of homeodomain transcription factors.

#### Single neuron-expressed genes

Between 160 (threshold 1, covering 44/128 neuron types) to 1,348 (threshold 4, covering 112/128 neuron types) genes are exclusively detected in a single neuron type (Table S3). The single-neuron specificities of many of these genes are validated by published, fosmid-based reporter gene analysis. For example, fosmid-based reporters for the *ceh-63* (DVA), *ceh-28* (M4), and *ceh-8* (RIA) homeobox genes match the neuron specificity of our scRNA-seq results (Reilly et al., 2020). The *cis*-regulatory control regions of these genes are candidate drivers for genetic access to individual cells in the nervous system (Lorenzo et al., 2020). Neurons not covered by single neuron-specific drivers can be genetically accessed by the intersection of drivers that are more broadly expressed.

#### Bulk RNA sequencing confirms scRNA-seq results and detects additional classes of non-coding RNAs

To validate our scRNA-seq dataset with an orthogonal approach, we used FACS to generate bulk RNA-seq profiles for eight neuron types: ASG, AVE, AVG, AWA, AWB, PVD, VD, and DD (Spencer et al., 2014; STAR Methods). Genes enriched in the single-cell clusters of these neurons (i.e., “marker genes”) were also most enriched in the corresponding bulk profiles (Figure 5A). For example, ASG marker genes from scRNA-seq (left column) are enriched ~24-fold ( $2^{4.61}$ ) in the ASG bulk RNA-seq profile (top left cell) compared to a pan-neuronal bulk reference. By contrast, markers for other cells are depleted in ASG bulk data (remainder of top row). Thus, independently derived single-cell and bulk RNA-seq datasets yielded consistent gene expression profiles. Consistent with their commingling in the scRNA-seq data, VD and DD GABAergic motor neurons had the fewest differentially expressed genes among all neuron pairs (Figure 5C). These results suggest that DD and VD GABAergic neurons are more closely related than are other pairs of different neuron types, and methods for distinguishing neuron types in single-cell data are relatively insensitive to small differences in gene expression.

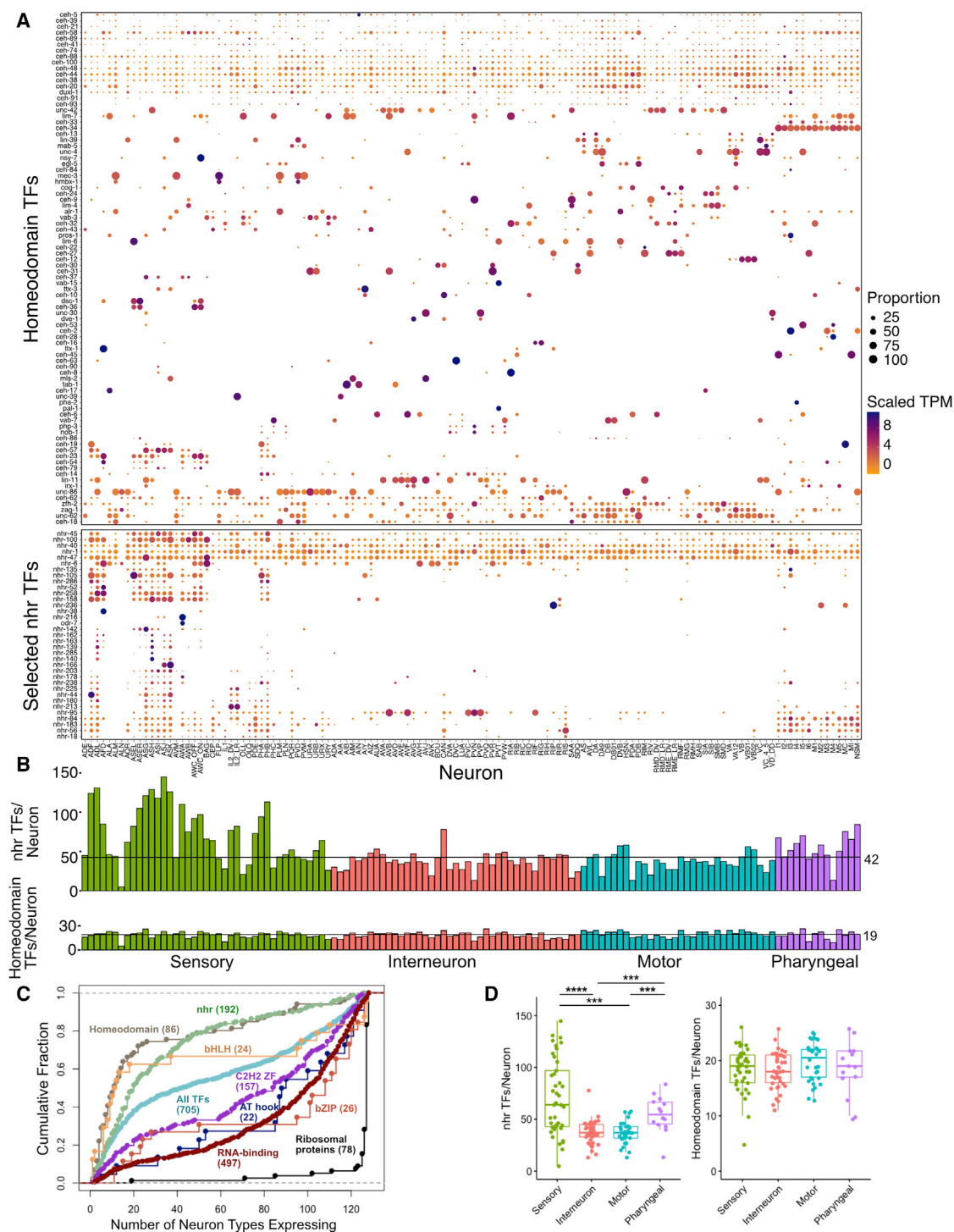
Protein coding genes, large intergenic noncoding RNAs (lincRNAs) and pseudogenes show similar coverage in both bulk and scRNA-seq datasets. However, as expected, non poly-adenylated non-coding RNAs (ncRNAs), small nuclear RNAs (snRNAs), and small nucleolar RNAs (snoRNAs) are rarely detected in our scRNA-seq data (possibly due to spurious priming) but are abundant in bulk RNA-seq samples derived from rRNA-depleted total RNA (Figure 5B). The smallest species of ncRNAs, microRNAs (miRNAs), and Piwi-interacting RNAs (piRNAs), are excluded from our bulk profiles due to a size exclusion step in library preparation, and their characterization awaits further studies.

#### Widespread differential splicing between neuron types

Differential splicing plays a critical role in the development and function of the nervous system (Raj and Blencowe, 2015; Vuong et al., 2016) and has been reported for individual neuron types in *C. elegans* (Moresco and Koelle, 2004; Norris et al., 2014;

(E) Number of all genes (top), neuropeptides (middle), and neuropeptide receptors (bottom) per neuron, grouped by neuron modality. Boxes are interquartile ranges. ANOVA, with Tukey post hoc comparisons for neuropeptide receptors, Kruskal-Wallis test for other comparisons. \* $p < 0.05$ , \*\*\* $p < 0.001$ . See also Figures S5 and S6, Tables S2 and S3, and Data S1.





**Figure 4. Expression of transcription factor families**

(A) Heatmap of homeodomain and representative subset of nuclear hormone receptor (nhr) transcription factors (TFs) across 128 neuron types (columns) grouped by neuron modality. TFs are clustered for each subfamily. Circle diameter represents the proportion of neurons in each cluster that expresses a given gene.

(B) Bar graphs of number of nhr and homeodomain TFs in each neuron type, grouped by neuron modality.

(C) Cumulative distribution of number of neuron types expressing homeodomain, bHLH, nhr, C2H2 ZF (zinc finger), AT hook, bZIP transcription factor (TF) families, RNA binding proteins, and ribosomal proteins (see also Figure 3A).

(legend continued on next page)

Thompson et al., 2019; Tomioka et al., 2016). Because the 3' bias of the 10x Genomics scRNA-seq method limits its use for detecting alternatively spliced transcripts (Arzalluz-Luque and Conesa, 2018; Dehghannasiri et al., 2020; Patrick et al., 2019), we leveraged the bulk RNA-seq profiles to identify differentially spliced transcripts among *C. elegans* neurons.

We discovered 111 high confidence occurrences of differential use of splicing sites between 8 neuron classes (Figures 5D–5F; Table S4). Most neuron pairs displayed some differential use of splicing sites (Figure 5D), with wide variations between pairs. For example, we detected 16 differential splicing events between ASG and VD, and only 2 differences between ASG and AWA.

In addition, we detected 63 previously unannotated exons (Table S4; STAR Methods). For example, the *mbk-2* transcript in AWA includes an additional 77-nt sequence corresponding to an alternative 5' exon that is not expressed in the other seven neuron types in our dataset (Figure 5F). This *mbk-2* exon is predicted by GenemarkHMM (Pavy et al., 1999), but its expression was not detected by whole-worm RNA-seq (Tourasse et al., 2017). Thus, our data underscore the capacity of bulk RNA-seq of single neuron types to detect differential splicing events that could not be reliably detected either by whole animal bulk RNA-seq or by 10x Genomics scRNA-seq.

### Analysis of *cis*-regulatory elements reveals a rich array of 5' and 3' motifs

To identify candidate *cis*-regulatory elements that underlie the distinct patterns of gene expression among neuron types, we used the FIRE motif discovery algorithm to analyze our scRNA-seq dataset. FIRE detects DNA motifs within promoter sequences and linear RNA motifs in 3' untranslated regions (UTRs) among cohorts of similarly regulated genes (Elemento et al., 2007). FIRE detects motifs that are significantly informative of relative gene expression in each neuron type (Figure 6A). Motifs of positive regulators, for example, should be significantly over-represented (yellow squares, red borders) in genes with high relative expression in the neuron (right columns). A subset of 5' DNA motifs matched known transcription factor DNA binding preferences (Khan et al., 2018; Weirauch et al., 2014). For example, a motif corresponding to the DNA binding sequence (CTACA) of several *nhr* transcription factors, including ODR-7, is over-represented in genes that are highly enriched in the AWA neuron (Figure 6A). Notably, ODR-7 is exclusively expressed in AWA where it regulates neuron identity (Colosimo et al., 2003; Sengupta et al., 1994, 1996).

We clustered all discovered motifs (see STAR Methods), resulting in 159 distinct DNA and 65 RNA motif families. 101 of 159 DNA motif families showed similarity to DNA binding sequences from available databases. For example, FIRE discovered a DNA motif family (TAATCC) which corresponds to the core DNA binding sequence of K50 class homeodomain transcription factors (Driever and Nüsslein-Volhard, 1989; Treisman et al., 1989) in genes with high relative expression in ASEL,

ASER, AWC<sup>ON</sup>, AWC<sup>OFF</sup>, BAG, and AWA neurons (Figure S7A). The TAATCC sequence matches *in vitro*-derived binding motifs for *C. elegans* K50 class homeodomain genes that are expressed in these neurons (*ceh-36* in ASE and AWC and *ceh-37* in BAG and AWA) (Figure S7A) and are required for their development (Chang et al., 2003; Koga and Ohshima, 2004; Lanjuin et al., 2003; Serrano-Saiz et al., 2013). These results indicate that our approach has the potential to reveal functionally relevant regulatory elements.

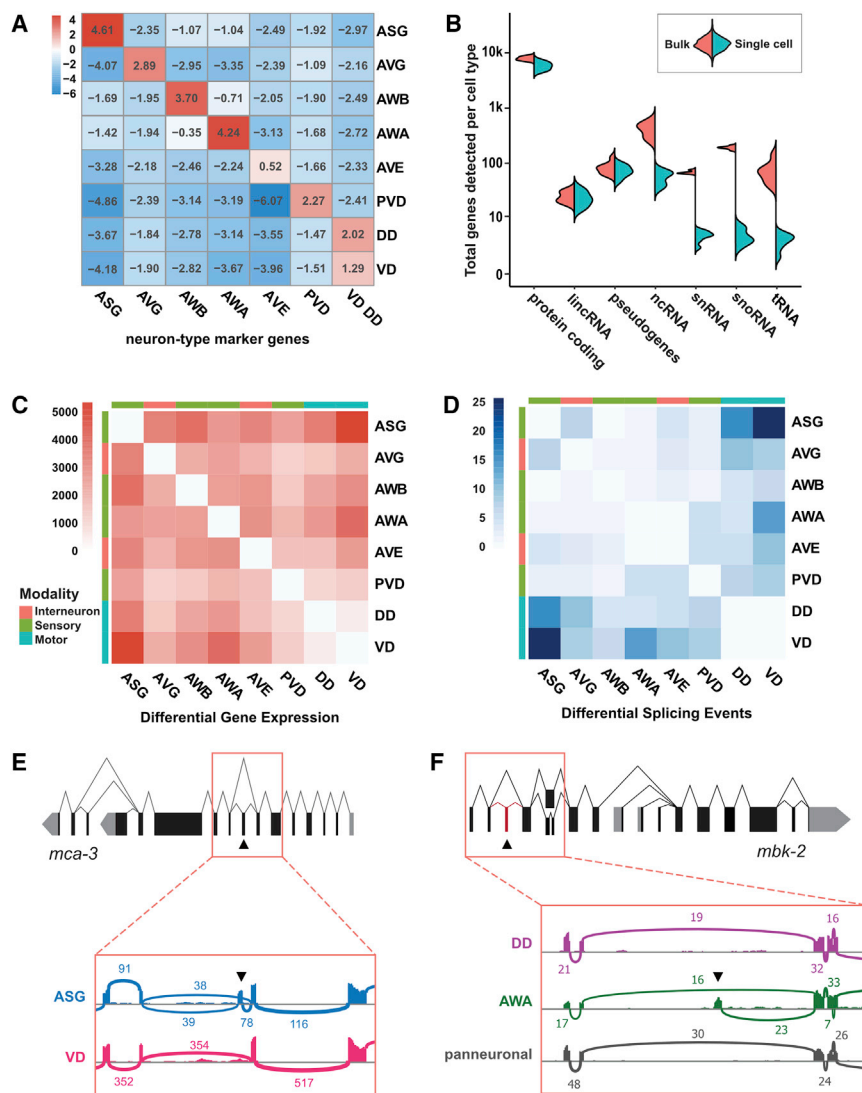
To limit false-positives, the FIRE algorithm uses stringent criteria for motif discovery and therefore generates conservative results. Although each motif family was discovered in an average of 5 neurons, we reasoned that the identified motif families might also regulate gene expression in additional neuron types. We therefore generated motif-neuron associations for each motif family (Figures 6B, 6C, and S7C; STAR Methods). We detected an average of 9 significant neuron associations for each motif family (log fold change >0.5 and p value <1e–5). This additional analysis significantly expanded the list of associations for neurons with previously established co-regulated genes. For example, motif family 184 matches the X-box sequences bound by DAF-19, which regulates cilia formation in all 28 ciliated neuron types (Efimenko et al., 2005; Swoboda et al., 2000). This X-box motif was initially discovered by FIRE in 10 ciliated neurons, but was significantly associated with another 12 ciliated sensory neurons by our additional analysis (Figure S7E).

Our approach also points to previously undetected roles for TFs in neuron-specific gene regulation. For example, motif family 85 corresponds to the E-box motif CAGGTG and is strongly associated with most amphid and phasmid neurons (Figure 6D). This particular E-box sequence is enriched in *hlh-4* target genes in the nociceptive sensory neuron ADL (Masoudi et al., 2018), but can also bind at least 10 distinct bHLH dimers (Grove et al., 2009). Interestingly, motif family 215 contained a different E-box sequence which was positively associated only with the chemorepulsive sensory neurons ADL, ASH, and PHB (Figure 6D). Based on the expression patterns of bHLH TFs in the adult nervous system, motif 215 may be a target of a HLH-2 homodimer (Masoudi et al., 2018).

Intriguingly, a substantial number of the motifs with strong positive associations with sensory neurons match TFs with uncharacterized roles in the nervous system or do not match any known TFs (Figure 6D). For example, motif family 100 showed a strong association with several sensory neurons and is similar to the binding site of the nuclear hormone receptor protein, NHR-142. *nhr-142* is almost exclusively expressed in a subset of amphid sensory neurons (Figure 4A), and the binding domain of *nhr-142* is closely related to several other *nhr* TFs (Lambert et al., 2019) that are expressed primarily in sensory neurons (*nhr-45*, *nhr-213*, *nhr-18*, *nhr-84*, and *nhr-178*), suggesting roles for these *nhr* TFs in sensory neuron function. Additionally, several motifs showed strong negative associations with enriched genes across many neurons (Figure 6D, right), indicating possible *cis*-regulatory elements of transcriptional repressors.

(D) Quantitative comparison of TFs per neuron for *nhr* (left) and homeodomain TFs (right) shows enrichment in sensory neurons for *nhrs*, but no differences for homeodomains. Boxplots are median and interquartile range (25<sup>th</sup>–75<sup>th</sup> percentile), Kruskal-Wallis. \*\*\*p < 0.001, \*\*\*\*p < 0.0001.

See also Figure S5, Tables S2 and S3, and Data S1.



**Figure 5. Comparison of bulk and single-cell RNA-seq**

(A) Heatmap for enrichment of scRNA-seq neuron-type marker genes (STAR Methods) (columns) in bulk RNA-seq data for each neuron type (ASG, AVG, AWB, AWA, AVE, PVD, DD, and VD) versus expression in all neurons. p values <0.001 for all comparisons except for AVE markers (all comparisons p value >0.05).

(B) Split violin plot quantifying detection of different RNA classes in bulk and scRNA-seq datasets for neuron types in (A).

(C and D) Heatmaps showing the number of differentially expressed genes (C) and differential splicing events (D) in pairwise comparisons of bulk RNA-seq datasets.

(E) Gene model and alternative splicing for *mca-3*. Inset: Sashimi plot shows alternative splicing of specific exon (arrowhead) in ASG versus VD.

(F) Gene model and alternative splicing of *mbk-2*. Inset: Sashimi plot shows detection of previously undescribed, alternatively spliced exon (arrowhead) in AWA but not in DD or pan neuronal bulk RNA-seq. For Sashimi plots in (E) and (F), vertical bars represent exonic reads and arcs indicate the number of junction-spanning reads. See also Table S4.

### Cell adhesion molecules are differentially expressed among neurons that are synaptically connected and that define anatomically distinct fascicles in the nerve ring

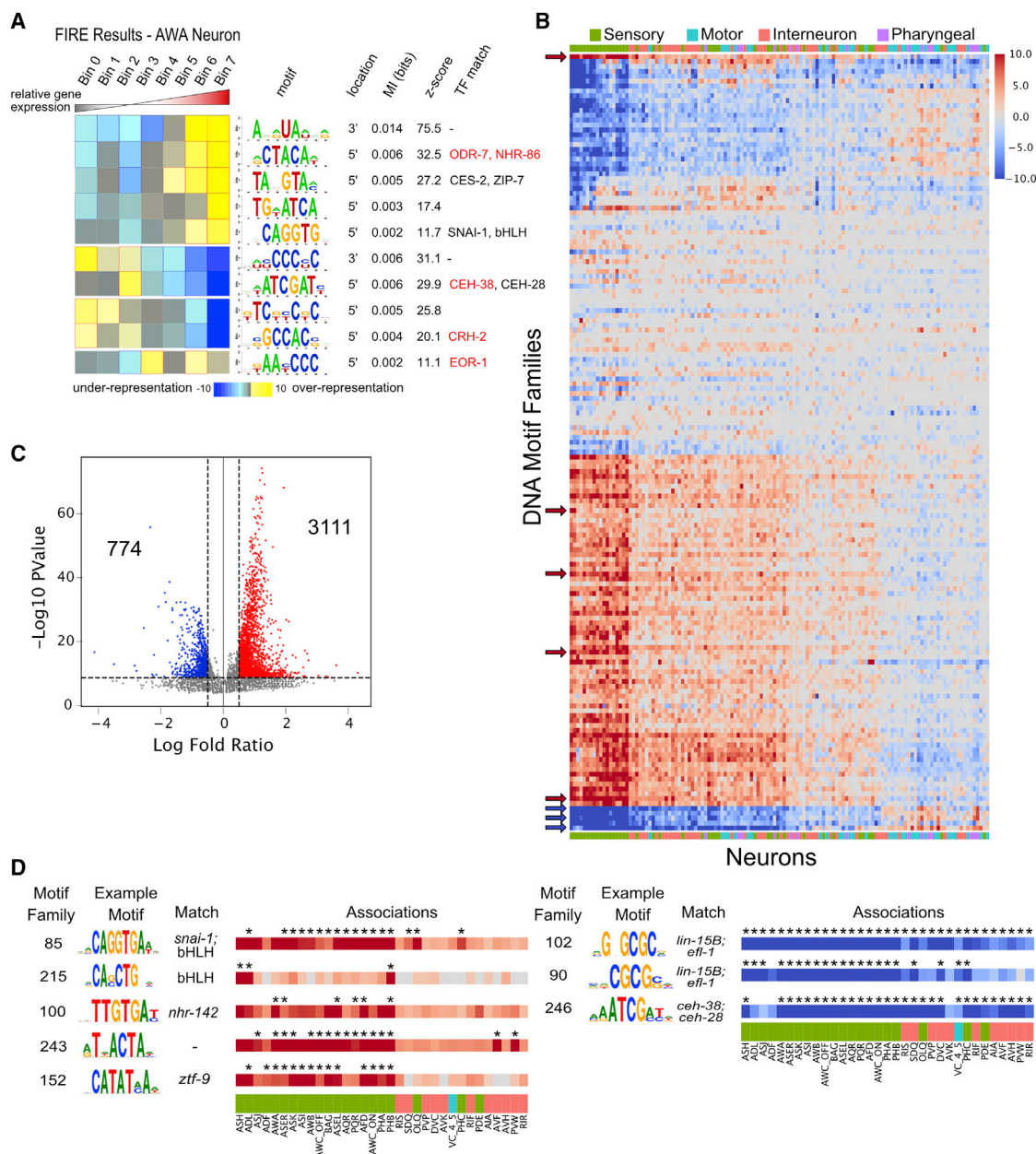
We compared our transcriptomic data to the *C. elegans* connectome to identify candidate genetic determinants of neurite bundling and synaptic connectivity. For this analysis, we utilized the nerve ring (Figure 7A), the largest expanse of neuropil in the *C. elegans* nervous system,

RNA motif analysis revealed that most RNA motif families showed positive associations with many neurons (indicating over-representation of RNA motifs in the enriched genes for each neuron type). Similar to DNA motifs, the strongest effects for RNA motifs were seen in sensory neurons (Figure S7F). In contrast to all other RNA motif families, motif family 23 showed negative associations with most neuron types. This motif family corresponds to a poly-C sequence (Figure S7G). A subclass of KH-domain RNA binding proteins interacts with poly-C regions in RNA and microRNAs (Choi et al., 2009). The *C. elegans* poly-C binding protein HRPK-1 positively regulates the function of several microRNA families, including those that act in the nervous system (Li et al., 2019). The over-representation of the poly-C motif family in depleted genes in most neurons indicates a potential role for this motif in microRNA-mediated repression. Overall, our analysis of neuron-specific gene expression identified over 200 *cis*-regulatory elements that could be sites for *trans*-acting factors such as transcription factors, RNA-binding proteins, and microRNAs.

because electron microscope reconstructions from multiple animals have detailed both membrane contacts and synapses in this region (Brittin et al., 2021; Cook et al., 2019; Witvliet et al., 2020). We limited our analyses to putative cell adhesion molecules (CAMs), which have documented roles in axon pathfinding, fasciculation, and synapse formation (Bruce et al., 2017; Colón-Ramos et al., 2007; Kim and Emmons, 2017; Shen and Bargmann, 2003; Siegenthaler et al., 2015; Sperry, 1963). 141 CAMs (Cox et al., 2004; Hobert, 2013; Table S3) were detected in neurons in our scRNA-seq dataset.

Recent computational analysis revealed a modular structure for the nerve ring, with four distinct neurite bundles or “strata” as well as a fifth group of unassigned neurons that contacts neurons in multiple strata (Moyle et al., 2021; see also Brittin et al., 2021; Figure S8A). Nerve ring formation begins in the embryo, but this structure is also modified throughout larval development as additional axons extend into the nerve ring and form synapses (Moyle et al., 2021; Witvliet et al., 2020). Together, these results





**Figure 6. Cis-regulatory elements in neuronal transcriptomes**

(A) FIRE results for AWA neuron, featuring the motif logo, location (5' or 3'), mutual information, z-scores from randomization-based statistical test and matching transcription factors. Genes were grouped into seven bins based on relative expression from lowest (left) to highest (right). Heatmap denotes over-representation (yellow) or under-representation (blue) of each motif (rows) in genes within each bin. Significant over-representation is indicated by red outlines, whereas significant under-representation is indicated by blue outlines. Transcription factors in red are expressed in AWA.

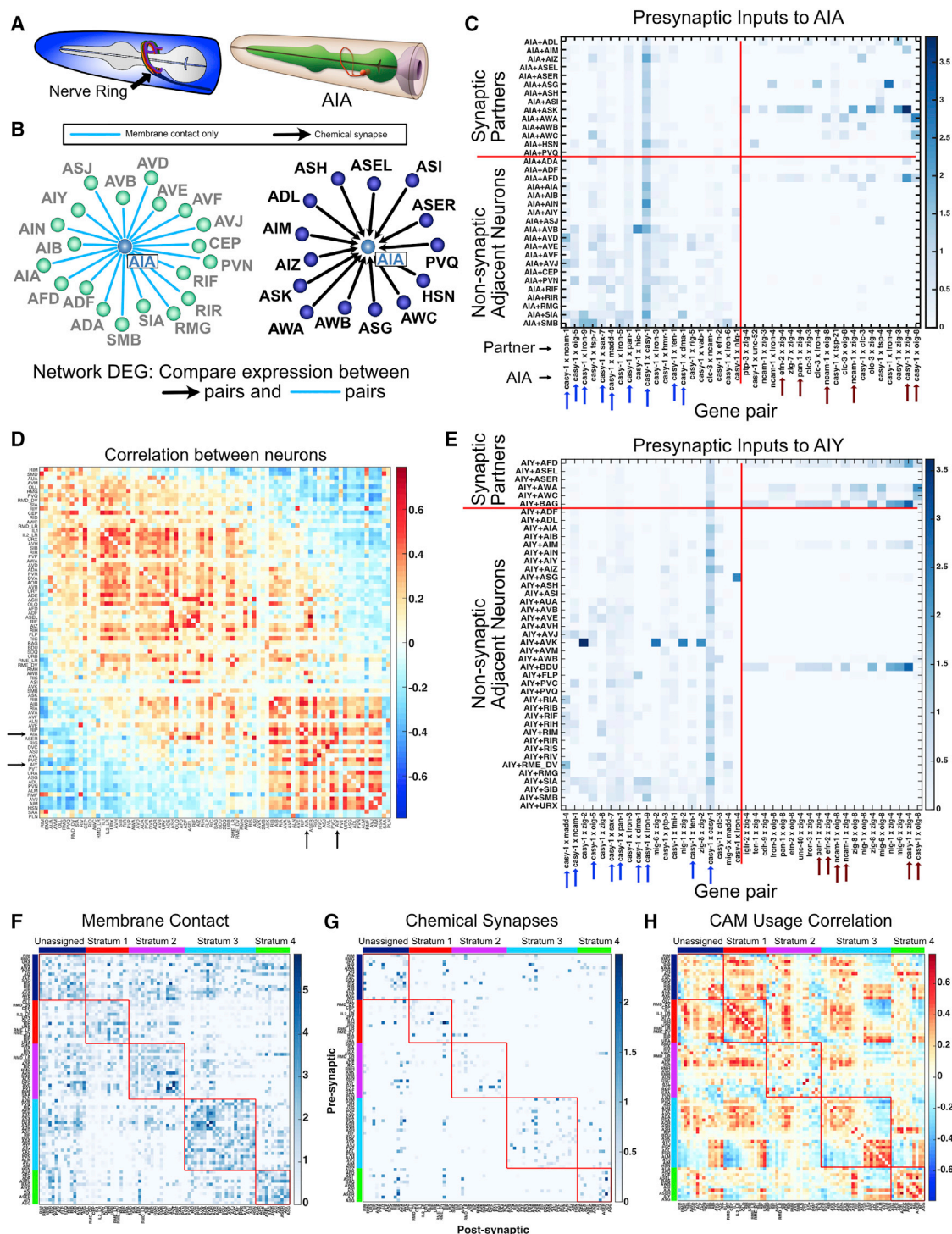
(B) Heatmap for enrichment of clustered motifs (rows) in each neuron class (columns). Red denotes enrichment in genes with highest relative expression, whereas blue indicates enrichment in genes with lowest relative expression (see STAR Methods). Color intensity represents  $\log_{10}(\text{p value})$  from hypergeometric test. Motif families and neurons are ordered by similarity. Color bar across x axis indicates neuron modality. Arrows denote motif families featured in (D).

(C) Volcano plot showing log fold ratio and  $-\log_{10} \text{p value}$  for all motif family-neuron associations. Significant associations with p value  $< 1e-5$  and log fold ratio  $> 0.5$  (3,111) or  $< -0.5$  (774) are noted.

(D) Eight selected motif families with significant associations with neurons from (C). Motif families: E-box motifs (85 and 215), motifs for nhrs (100), homeodomains (246), and a previously undescribed motif (243). Asterisks denote significant associations.

See also Figure S7.





**Figure 7. Differential expression of cell adhesion molecules among neurons and their presynaptic partners**

(A) Left: the *C. elegans* nerve ring. Right: AIA ring interneuron. From WormAtlas.

(B) Neurons with presynaptic input to AIA (right) and neurons with membrane contact but no synapses with AIA (left).

(C) Heatmap of 20 cell adhesion molecule (CAM) gene pairs with highest log fold change in AIA + presynaptic inputs versus AIA + non-synaptic adjacent neurons (right of vertical red line). 20 CAM gene pairs with highest log fold change in AIA + non-synaptic adjacent neurons versus AIA + presynaptic partners (left of vertical red line). Arrows denote gene pairs common for AIA and AIY (E).

(D) Correlation matrix for CAM usage (see text) across all neurons in the nerve ring (84 neuron types). Arrows indicate AIA and AIY (correlation = 0.568).

(E) Heatmap as in (C), for AIY. Arrows denote gene pairs common for AIA and AIY.

(legend continued on next page)

point to the importance of both periodic as well as sustained expression of genetic determinants that initiate, modify or maintain the overall structure of the nerve ring and its connectome.

We first determined CAMs that were differentially expressed between strata (Figures S8B and S8C). Six CAMs were significantly enriched in the neurons in one stratum compared to the neurons in all other strata (Figure S8C). Notably, the transcript for MADD-4/punctin, a secreted protein that has been shown to direct process outgrowth as well synaptic placement (Zhou and Bessereau, 2019), is significantly enriched in stratum 1. *tsp-7*, a homolog of the human protein CD63, a member of the tetraspanin superfamily, is highly expressed in stratum 2. Tetraspanins interact with integrins and have been implicated in membrane trafficking and synaptogenesis (Murru et al., 2018; Pols and Klumperman, 2009). *Iron-5* and *Iron-9* (extracellular leucine rich repeat proteins) are selectively expressed in a subset of neurons in stratum 2 that could be indicative of roles in organizing these specific fascicles (Figure S8B). Thus, our approach has identified candidate genes that can now be experimentally tested for roles in organizing and maintaining structurally and functionally distinct domains of the nerve ring.

In addition to mediating axon fasciculation, we reasoned that specific CAMs might contribute to synaptic maintenance in the mature nervous system. We surmised that CAMs mediating synaptic stability are more highly expressed in synaptically connected neurons than in adjacent neurons with membrane contacts but no synapses. We generated high-confidence membrane adjacency and chemical synaptic connectomes by retaining only contacts and synapses that are preserved across animals in EM reconstructions of the nerve ring (Table S5; STAR Methods; Brittin et al., 2021; Cook et al., 2019; White et al., 1986; Witvliet et al., 2020). These datasets include 84 of the 128 neuron classes. The importance of genetic determinants of connectivity in this circuit is underscored by the observation that membrane contacts between neurons in the nerve ring are much more numerous than synapses; on average, in the nerve ring, each neuron synapses with only 15% of the neurons it contacts (means of 6.42 presynaptic inputs, 6.42 postsynaptic outputs, 42 contacted cells) (Brittin et al., 2021; White et al., 1986).

For each neuron, we compared the expression of all possible combinations of pairs of CAMs in the neuron and its synaptic partners relative to the neuron and its non-synaptic adjacent neurons (Figures 7B and 7C). Two independent comparisons were generated, one for presynaptic partners (Figure 7C) and a second result for postsynaptic neurons (Methods S1). Our analysis revealed multiple CAM gene pairs with enrichment in synaptically connected neurons compared to adjacent but not synaptically connected neurons. A representative example for presynaptic inputs to the interneuron AIA shows that CAM pairs enriched in synaptically connected neurons were not uniform for the different presynaptic partners of AIA (Figure 7C). For example, AIA and its presynaptic partner, ASK, show strong enrichment for *casy-1* (calsyntenin)

and *zig-4* (secreted 2-immunoglobulin [Ig] domain protein) whereas the AIA-ASG pair is enriched for *casy-1* (calsyntenin) and *Iron-4* (extracellular leucine rich repeat protein). This finding is consistent with the prediction that distinct combinatorial codes of CAMs could be required for patterning connectivity between individual pairs of neurons (Kim and Emmons, 2017). Additionally, we identified distinct CAM pairs that are enriched in adjacent, not synaptically connected neurons (Figure 7C). This observation indicates that some CAM interactions may functionally inhibit either the formation or maintenance of synapses between neurons. Anti-synaptic effects have been documented for the axon guidance molecules netrin, sema-5B, and their cell surface receptors (O'Connor et al., 2009; Poon et al., 2008; Tran et al., 2009).

To examine patterns across the nerve ring, we restricted our analysis to gene pairs with a log fold change >0.2 in either synaptically connected or in adjacent but not connected neurons for at least one neuron type. We refer to this pattern of CAM pairs enriched in synaptic or solely adjacent neurons as “CAM usage.” Of 19,881 possible CAM pairs, 439 pairs passed our log fold change threshold for presynaptic connections, whereas 443 pairs showed >0.2 log fold change for postsynaptic connections (Methods S1). To identify neurons with similar patterns of presynaptic CAM usage, we generated correlation matrices from pairwise comparisons of all neurons and sorted neurons by similarity using multidimensional scaling (Figure 7D). For example, CAM usage for presynaptic inputs to AIA and AIY is strongly correlated (correlation, 0.568) due to the co-occurrence for each neuron of multiple shared combinations of CAMs (Figure 7E, blue and red arrows). This analysis also separated neurons into two main groups based on CAM usage that could be indicative of underlying shared roles for CAMs among these distinct sets of neurons.

We sought to understand the relationship between stratum membership and synaptic CAM usage for nerve ring neurons. Both membrane contact and chemical synapses are denser among neurons within strata than across strata (Figures 7F and 7G), a finding also observed for an independent assessment of nerve ring axon bundles (Brittin et al., 2021). We sorted neurons by CAM usage within each stratum (Figure 7H) to assess intra-stratum correlations. This approach revealed high correlations among neurons within strata. Additionally, neurons in some strata split into distinct groups based on CAM usage (Stratum 3) (Figure 7H; Methods S1). This observation suggests that CAM usage at synaptic connections is likely distinct from CAMs that may be involved in strata formation and/or maintenance. Although CAM usage correlations were often elevated among neurons within strata, high correlations were also detected among neurons in different strata that are not synaptically connected and with minimal contacts, thus suggesting roles for CAMs in nerve ring architecture and connectivity likely depend on additional factors. We suggest that the overall results of our analysis point to specific CAMs that can now be investigated for roles in the formation

(F) Membrane adjacency matrix was grouped by nerve ring strata (each outlined with red box) (Moyle et al., 2021). Within each stratum, neurons were ordered according to CAM usage correlations (see H).

(G) Strata ordering as in (F) was imposed upon the chemical connectome revealing that most synapses are detected between neurons within the same stratum.

(H) The CAM usage correlation matrix (as in D) was grouped by strata, then sorted by similarity within each stratum. CAM usage is broadly shared for neurons in strata 1 and 4. Stratum 3 shows two distinct populations.

See also Figure S8, Table S5, and Methods S1.

and maintenance of synapses as well as fasciculation between specific neurons in the *C. elegans* nerve ring.

### Data interface

We developed a web application, CengenApp (<https://cengen.shinyapps.io/CengenApp>) to facilitate analysis of these scRNA-seq data. Users can generate gene expression profiles by neuron class or by gene at different thresholds, and perform differential gene expression analysis between either individual neurons or between groups of neuron types. In addition, an interactive graphical interface is available for generating heatmap representations (e.g., Figure 3C) of gene expression across the nervous system. Raw data are available at Gene Expression Omnibus (<https://www.ncbi.nlm.nih.gov/geo>) (single cell data at GEO: GSE136049, bulk data at GEO: GSE169137). The data and additional supporting files can be downloaded from the CeNGEN website (<https://www.cengen.org>) and code is available at GitHub (<https://www.github.com/cengenproject>).

### Conclusions

We have produced a gene expression map for the entire *C. elegans* nervous system, complementing earlier partial profiles of the *C. elegans* nervous system at embryonic and early larval stages (Cao et al., 2017; Packer et al., 2019). This catalog of gene expression provides an essential foundation for a comprehensive exploration of transcriptional and gene regulatory patterns that lead to neuronal diversity, connectivity, and function. *C. elegans* is the first organism in which a complete anatomical map of its nervous system is matched with a nervous system-wide molecular map, therefore providing new opportunities to investigate neuronal development and function.

We developed a thresholding approach for single-cell data to generate high confidence profiles for each neuron type. Multiple findings indicate that neuropeptide signaling is widely utilized and likely crucial for a variety of functions. First, neuropeptide-encoding genes are among the most abundantly detected genes in the dataset. Second, at the most stringent threshold examined, each neuron expresses at least four different neuropeptide-encoding genes. Third, each neuron expresses a distinct combination of both neuropeptide genes and putative neuropeptide receptors. Recent reports show abundant and widespread neuropeptide expression in *Hydra* (Siebert et al., 2019), *Drosophila* (Allen et al., 2020), and mouse cortical neurons (Smith et al., 2019), indicating that these salient features of neuropeptide signaling are conserved among diverse species.

Our analysis of transcription factor expression reveals that different transcription factor families appear to have segregated into distinct functions during cellular differentiation. Some families are underrepresented in the mature nervous system (T-box genes), others show broad expression patterns in the nervous system (Zn finger), whereas others are sparsely expressed and appear to exquisitely track with neuronal identity (homeodomain) (Reilly et al., 2020). The nuclear hormone receptors (nhrs) may have acquired a unique function, as inferred by their striking enrichment in sensory neurons. The identification of enriched *cis*-regulatory motifs in neuronal gene batteries provides an opportunity for future experiments to dissect the mechanisms of gene regulation in the nervous system.

Finally, we devised computational strategies that exploit our gene expression profile of the *C. elegans* nervous system to reveal the genetic underpinnings of neuron-specific process placement and connectivity. Previous computational efforts to forge a link between neuron-specific gene expression and the *C. elegans* wiring diagram have been hampered by incomplete and largely qualitative expression data (Barabási and Barabási, 2020; Baruch et al., 2008; Kaufman et al., 2006; Kovacs et al., 2020; Varadan et al., 2006). Here, we leveraged our nervous-system wide catalog of gene expression to deduce combinatorial codes for CAMs that likely contribute to the maintenance and formation of this complex neuropil. Importantly, this analysis can now be extended to specific groups of neurons and to any gene family to generate specific hypotheses of process placement and connectivity for direct experimental validation.

We expect that these data will be useful for future studies of individual genes, neurons, and circuits, as well as global analyses of an entire nervous system and the development of scRNA-seq analysis methods. Coupled with the fully described cell lineages (Sulston and Horvitz, 1977; Sulston et al., 1983), neuronal anatomy (Albertson and Thomson, 1976; Brittin et al., 2021; Cook et al., 2019; White et al., 1986; Witvliet et al., 2020), and powerful functional analyses, such as pan-neuronal calcium imaging and neuronal identification (Kato et al., 2015; Nguyen et al., 2016; Venkatachalam et al., 2016; Yemini et al., 2021), our dataset provides the foundation for discovering the genetic programs underlying neuronal development, connectivity, and function across an entire nervous system.

### Limitations of the study

Although we provide gene expression profiles of every neuron class in the *C. elegans* hermaphrodite, these neuron-specific transcriptomes are incomplete for several reasons:

- (1) Some neuron classes are under-represented, likely due to biases in the dissociation procedure, thus resulting in incomplete detection of expressed transcripts in the corresponding scRNA-seq dataset (Figures S5I–S5L).
- (2) Our scRNA-seq library construction method largely excluded non-coding RNAs that are not poly-adenylated (Figure 5B).
- (3) Alternative splicing is rarely detected in our scRNA-seq dataset due to short reads and the 3' bias of the library construction method (Figures 5D–5F).

Additional approaches, such as isolation of individual neuron types for bulk RNA-seq (Figure 5A), single-nuclei RNA-seq, long-read sequencing and alternative RNA-seq library preparation methods could be used in future studies to produce a more comprehensive description of the *C. elegans* neuronal transcriptome.

### STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE



- **RESOURCE AVAILABILITY**
  - Lead contact
  - Materials availability
  - Data and code availability
- **EXPERIMENTAL MODEL AND SUBJECT DETAILS**
  - Preparation of larvae and dissociation
- **METHOD DETAILS**
  - FACS isolation of neuron types for RNA-Seq
  - Single-cell RNA sequencing
  - Single-cell RNA-Seq Mapping
  - Downstream Processing
  - Dimensionality reduction and batch correction
  - Cell Identification
  - Neuron network analysis
  - Gene expression analyses
  - Stress-induced genes
  - Thresholding
  - Estimating coverage for individual neurons
  - Determining distinct combinations of gene sets
  - Connectivity Analysis
  - Reporter strains
  - Imaging
  - RNA Extraction
  - Bulk sequencing and mapping
  - Comparing scRNA-Seq and bulk RNA data
  - Alternative splicing
  - Generating connectivity matrices
  - Regulatory patterns of neuron transcriptomes
  - Cis-regulatory element discovery
  - Motif families
  - Associations of motif families and neurons
  - Cell adhesion molecule by stratum analysis
  - Network differential gene expression analysis
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
- **ADDITIONAL RESOURCES**

## SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.cell.2021.06.023>.

## ACKNOWLEDGMENTS

We thank the CeNGEN Advisory Board for guidance, M. Zhen for ZM9592, and H. Sun for imaging neuropeptide reporters. O.H. is an Investigator with the Howard Hughes Medical Institute. FACS (Flow Cytometry Shared Resource, supported by Ingram Cancer Center [P30 CA68485], DDRRC [DK058404]), scRNA-seq (VANTAGE, supported by CTSA [5UL1 RR024975-03], Ingram Cancer Center [P30 CA68485], Vision Center [P30 EY08126], and NIH/NCRR [G20 RR030956]), and confocal imaging (Cell Imaging Shared Resource [NIH CA68485, DL20593, DK58404, DK59637, and EY08126]) were performed at Vanderbilt. Strains were provided by the CGC (NIH P40 OD010440A). G.S. was supported by “la Caixa” Foundation (LCF/BQ/PI19/11690010, ID 100010434) and by Ministerio de Ciencia e Innovación, Spain (PID2019-104700GA-I00). This work was funded by NIH (R01NS100547 to M.H., O.H., D.M.M., and N.S. and R01 NS110391 to O.H.) and by Vanderbilt TIPs (to D.M.M.).

## AUTHOR CONTRIBUTIONS

M.H., O.H., N.S., and D.M.M. originated project. S.R.T. generated scRNA-seq data; assigned neuron identities; analyzed gene family expression; helped

A.W., C.X., E.V., and P.O. with data analysis; designed figures and tables; wrote the first draft; and edited the final version. G.S. developed CengenApp. A.W. designed thresholding strategy with S.R.T., analyzed alternative splicing, and implemented meta data format. A.B. generated and analyzed bulk RNA sequence data. M.B.R. provided ground truth reporters. C.X. extended 3' UTRs for read mapping. E.V. correlated CAMs with neuron-specific synapses and strata and provided input on statistical and quantitative analysis. P.O. implemented FIRE analysis with S.T. L.G. provided BrainAtlas. R.M. and A.P. used FACS to isolate neurons. R.M. extracted RNA for bulk RNA-seq. I.R., E.Y., S.J.C., B.V., C.C., M.B., A.A., and S.R.T. generated reporter strains. E.Y. helped with NeuroPAL. N.S. directed G.S. and C.X. and edited the manuscript. M.H. directed G.S., A.W., A.B., M.B., and A.A.; contributed to the first draft; and edited the final version. O.H. directed M.B.R., L.G., I.R., E.Y., S.C., B.V., and C.C.; contributed to the first draft; and edited the final version. D.M.M. oversaw work; directed S.R.T., R.M., and A.P.; contributed to the first draft; and edited the final version.

## DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: December 14, 2020

Revised: April 9, 2021

Accepted: June 14, 2021

Published: July 7, 2021

## WEB RESOURCES

CeNGEN, <https://www.cengen.org>

CenGEN Github, <https://www.github.com/cengenproject>

CengenApp, <https://cengen.shinyapps.io/CengenApp>

The Peptide-GPCR Project, <https://worm.peptide-gpcr.org/project>

## SUPPORTING CITATIONS

The following references appear in the supplemental information: [Dlakić \(2002\)](#); [Finak et al. \(2015\)](#); [Robinson et al. \(2010\)](#); [Wang et al. \(2019\)](#).

## REFERENCES

- Adorjan, I., Tyler, T., Bhaduri, A., Demharter, S., Finszter, C.K., Bako, M., Se-bok, O.M., Nowakowski, T.J., Khodosevich, K., Möllgård, K., et al. (2019). Neuroserpin expression during human brain development and in adult brain revealed by immunohistochemistry and single cell RNA sequencing. *J. Anat.* 235, 543–554.
- Albertson, D.G., and Thomson, J.N. (1976). The pharynx of *Caenorhabditis elegans*. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 275, 299–325.
- Allen, A.M., Neville, M.C., Birtles, S., Croset, V., Treiber, C.D., Waddell, S., and Goodwin, S.F. (2020). A single-cell transcriptomic atlas of the adult *Drosophila* ventral nerve cord. *eLife* 9, e54074.
- Arzalluz-Luque, Á., and Conesa, A. (2018). Single-cell RNAseq for the study of isoforms-how is that possible? *Genome Biol.* 19, 110.
- Barabási, D.L., and Barabási, A.L. (2020). A Genetic Model of the Connectome. *Neuron* 105, 435–445.e5.
- Barrett, A., McWhirter, R., Taylor, S.R., Weinreb, A., Miller, III, D.M., and Hammarlund, M. (2021). A head-to-head comparison of ribodepletion and polyA selection approaches for *C. elegans* low input RNA-sequencing libraries. *G3*. <https://doi.org/10.1093/g3journal/jkab121>.
- Baruch, L., Itzkovitz, S., Golan-Mashiach, M., Shapiro, E., and Segal, E. (2008). Using expression profiles of *Caenorhabditis elegans* neurons to identify genes that mediate synaptic connectivity. *PLoS Comput. Biol.* 4, e1000120.
- Becht, E., McInnes, L., Healy, J., Dutertre, C.A., Kwok, I.W.H., Ng, L.G., Ginhoux, F., and Newell, E.W. (2018). Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* 37, 38–44.



- Bhattacharya, A., Aghayeva, U., Berghoff, E.G., and Hobert, O. (2019). Plasticity of the Electrical Connectome of *C. elegans*. *Cell* 176, 1174–1189.e16.
- Brenner, S. (1974). The genetics of *Caenorhabditis elegans*. *Genetics* 77, 71–94.
- Brittin, C.A., Cook, S.J., Hall, D.H., Emmons, S.W., and Cohen, N. (2021). A multi-scale brain map derived from whole-brain volumetric reconstructions. *Nature* 591, 105–110.
- Bruce, F.M., Brown, S., Smith, J.N., Fuerst, P.G., and Erskine, L. (2017). DSCAM promotes axon fasciculation and growth in the developing optic pathway. *Proc. Natl. Acad. Sci. USA* 114, 1702–1707.
- Brunquell, J., Morris, S., Lu, Y., Cheng, F., and Westerheide, S.D. (2016). The genome-wide role of HSF-1 in the regulation of gene expression in *Caenorhabditis elegans*. *BMC Genomics* 17, 559.
- Canty, A., and Ripley, B. (2019). boot: Bootstrap R (S-Plus) Functions. R package version 1.3-24.
- Cao, J., Packer, J.S., Ramani, V., Cusanovich, D.A., Huynh, C., Daza, R., Qiu, X., Lee, C., Furlan, S.N., Steemers, F.J., et al. (2017). Comprehensive single-cell transcriptional profiling of a multicellular organism. *Science* 357, 661–667.
- Cao, J., Spielmann, M., Qiu, X., Huang, X., Ibrahim, D.M., Hill, A.J., Zhang, F., Mundlos, S., Christiansen, L., Steemers, F.J., et al. (2019). The single-cell transcriptional landscape of mammalian organogenesis. *Nature* 566, 496–502.
- Chang, S., Johnston, R.J., Jr., and Hobert, O. (2003). A transcriptional regulatory cascade that controls left/right asymmetry in chemosensory neurons of *C. elegans*. *Genes Dev.* 17, 2123–2137.
- Choi, H.S., Hwang, C.K., Song, K.Y., Law, P.Y., Wei, L.N., and Loh, H.H. (2009). Poly(C)-binding proteins as transcriptional regulators of gene expression. *Biochem. Biophys. Res. Commun.* 380, 431–436.
- Colón-Ramos, D.A., Margeta, M.A., and Shen, K. (2007). Glia promote local synaptogenesis through UNC-6 (netrin) signaling in *C. elegans*. *Science* 318, 103–106.
- Colosimo, M.E., Tran, S., and Sengupta, P. (2003). The divergent orphan nuclear receptor ODR-7 regulates olfactory neuron gene expression via multiple mechanisms in *Caenorhabditis elegans*. *Genetics* 165, 1779–1791.
- Cook, S.J., Jarrell, T.A., Brittin, C.A., Wang, Y., Bloniarz, A.E., Yakovlev, M.A., Nguyen, K.C.Q., Tang, L.T.H., Bayer, E.A., Duerr, J.S., et al. (2019). Whole-animal connectomes of both *Caenorhabditis elegans* sexes. *Nature* 571, 63–71.
- Cox, E.A., Tuskey, C., and Hardin, J. (2004). Cell adhesion receptors in *C. elegans*. *J. Cell Sci.* 117, 1867–1870.
- Csardi, G., and Nepusz, T. (2006). The igraph software package for complex network research. *InterJournal Complex Systems*, 1695.
- Davie, K., Janssens, J., Koldere, D., De Waegeneer, M., Pech, U., Kreft, L., Aibar, S., Makhzami, S., Christiaens, V., Bravo González-Blas, C., et al. (2018). A Single-Cell Transcriptome Atlas of the Aging *Drosophila* Brain. *Cell* 174, 982–998.e20.
- Davison, A.C., and Hinkley, D.V. (1997). *Bootstrap Methods and Their Applications* (Cambridge University Press).
- Dehghannasiri, R., Olivieri, J.E., and Salzman, J. (2020). Specific splice junction detection in single cells with SICILIAN. *bioRxiv*. <https://doi.org/10.1101/2020.04.14.041905>.
- Diakić, M. (2002). A new family of putative insulin receptor-like proteins in *C. elegans*. *Curr. Biol.* 12, R155–R157.
- Driever, W., and Nüsslein-Volhard, C. (1989). The bicoid protein is a positive regulator of hunchback transcription in the early *Drosophila* embryo. *Nature* 337, 138–143.
- Efimenko, E., Bubb, K., Mak, H.Y., Holzman, T., Leroux, M.R., Ruvkun, G., Thomas, J.H., and Swoboda, P. (2005). Analysis of *xbx* genes in *C. elegans*. *Development* 132, 1923–1934.
- Elemento, O., Slonim, N., and Tavazoie, S. (2007). A universal framework for regulatory element discovery across all genomes and data types. *Mol. Cell* 28, 337–350.
- Finak, G., McDavid, A., Yajima, M., Deng, J., Gersuk, V., Shalek, A.K., Slichter, C.K., Miller, H.W., McElrath, M.J., Pric, M., et al. (2015). MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol.* 16, 278.
- Fuxman Bass, J.I., Pons, C., Kozłowski, L., Reece-Hoyes, J.S., Shrestha, S., Holdorf, A.D., Mori, A., Myers, C.L., and Walhout, A.J. (2016). A gene-centered *C. elegans* protein-DNA interaction network provides a framework for functional predictions. *Mol. Syst. Biol.* 12, 884.
- Gendrel, M., Atlas, E.G., and Hobert, O. (2016). A cellular and regulatory map of the GABAergic nervous system of *C. elegans*. *eLife* 5, e17686.
- 10X Genomics (2017). Transcriptional profiling of 1.3 million brain cells with the Chromium Single Cell Gene Expression Solution. [https://pages.10xgenomics.com/rs/446-PBO-704/images/10x\\_LIT015\\_Chromium\\_Million-Brain-Cells\\_Application-Note\\_Letter\\_digital.pdf](https://pages.10xgenomics.com/rs/446-PBO-704/images/10x_LIT015_Chromium_Million-Brain-Cells_Application-Note_Letter_digital.pdf).
- Granato, M., Schnabel, H., and Schnabel, R. (1994). pha-1, a selectable marker for gene transfer in *C. elegans*. *Nucleic Acids Res.* 22, 1762–1763.
- Grove, C.A., De Masi, F., Barrasa, M.I., Newburger, D.E., Alkema, M.J., Bulky, M.L., and Walhout, A.J.M. (2009). A multiparameter network reveals extensive divergence between *C. elegans* bHLH transcription factors. *Cell* 138, 314–327.
- Hafemeister, C., and Satija, R. (2019). Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.* 20, 296.
- Haghverdi, L., Lun, A.T.L., Morgan, M.D., and Marioni, J.C. (2018). Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* 36, 421–427.
- Hallam, S., Singer, E., Waring, D., and Jin, Y. (2000). The *C. elegans* NeuroD homolog *cnd-1* functions in multiple aspects of motor neuron fate specification. *Development* 127, 4239–4252.
- Hammarlund, M., Hobert, O., Miller, D.M., 3rd, and Sestan, N. (2018). The CeNGEN Project: The Complete Gene Expression Map of an Entire Nervous System. *Neuron* 99, 430–433.
- Harris, T.W., Arnaboldi, V., Cain, S., Chan, J., Chen, W.J., Cho, J., Davis, P., Gao, S., Grove, C.A., Kishore, R., et al. (2020). WormBase: a modern Model Organism Information Resource. *Nucleic Acids Res.* 48 (D1), D762–D767.
- Hirose, T., Galvin, B.D., and Horvitz, H.R. (2010). Six and Eya promote apoptosis through direct transcriptional activation of the proapoptotic BH3-only gene *egl-1* in *Caenorhabditis elegans*. *Proc. Natl. Acad. Sci. USA* 107, 15479–15484.
- Hobert, O. (2002). PCR fusion-based approach to create reporter gene constructs for expression analysis in transgenic *C. elegans*. *Biotechniques* 32, 728–730.
- Hobert, O. (2013). The neuronal genome of *Caenorhabditis elegans*. *WormBook*, 1–106.
- Hobert, O., Glenwinkel, L., and White, J. (2016). Revisiting Neuronal Cell Type Classification in *Caenorhabditis elegans*. *Curr. Biol.* 26, R1197–R1203.
- Husson, S.J., Lindemans, M., Janssen, T., and Schoofs, L. (2009). Comparison of *Caenorhabditis elegans* NLP peptides with arthropod neuropeptides. *Trends Parasitol.* 25, 171–181.
- Inoue, T., Sherwood, D.R., Aspöck, G., Butler, J.A., Gupta, B.P., Kirouac, M., Wang, M., Lee, P.Y., Kramer, J.M., Hope, I., et al. (2002). Gene expression markers for *Caenorhabditis elegans* vulval cells. *Mech. Dev.* 119 (Suppl 1), S203–S209.
- Johnston, R.J., Jr., Chang, S., Etchberger, J.F., Ortiz, C.O., and Hobert, O. (2005). MicroRNAs acting in a double-negative feedback loop to control a neuronal cell fate decision. *Proc. Natl. Acad. Sci. USA* 102, 12449–12454.
- Kahles, A., Ong, C.S., Zhong, Y., and Ratsch, G. (2016). SplAdder: identification, quantification and testing of alternative splicing events from RNA-Seq data. *Bioinformatics* 32, 1840–1847.
- Kaletsky, R., Lakhina, V., Arey, R., Williams, A., Landis, J., Ashraf, J., and Murphy, C.T. (2016). The *C. elegans* adult neuronal IIS/FOXO transcriptome reveals adult phenotype regulators. *Nature* 529, 92–96.
- Kato, S., Kaplan, H.S., Schrödel, T., Skora, S., Lindsay, T.H., Yemini, E., Lockery, S., and Zimmer, M. (2015). Global brain dynamics embed the motor command sequence of *Caenorhabditis elegans*. *Cell* 163, 656–669.

- Kaufman, A., Dror, G., Meilijson, I., and Rupp, E. (2006). Gene expression of *Caenorhabditis elegans* neurons carries information on their synaptic connectivity. *PLoS Comput. Biol.* 2, e167.
- Kerk, S.Y., Kratsios, P., Hart, M., Mourao, R., and Hobert, O. (2017). Diversification of *C. elegans* Motor Neuron Identity via Selective Effector Gene Repression. *Neuron* 93, 80–98.
- Khan, A., Fornes, O., Stigliani, A., Gheorghe, M., Castro-Mondragon, J.A., van der Lee, R., Bessy, A., Chèneby, J., Kulkarni, S.R., Tan, G., et al. (2018). JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Res.* 46 (D1), D260–D266.
- Kim, B., and Emmons, S.W. (2017). Multiple conserved cell adhesion protein interactions mediate neural wiring of a sensory circuit in *C. elegans*. *eLife* 6, e29257.
- Koga, M., and Ohshima, Y. (2004). The *C. elegans* *ceh-36* gene encodes a putative homeodomain transcription factor involved in chemosensory functions of ASE and AWC neurons. *J. Mol. Biol.* 336, 579–587.
- Kovacs, I.A., Barabási, D.L., and Barabási, A.L. (2020). Uncovering the genetic blueprint of the *C. elegans* nervous system. *bioRxiv*. <https://doi.org/10.1101/2020.05.04.076315>.
- Kozioł, U., Kozioł, M., Preza, M., Costabile, A., Brehm, K., and Castillo, E. (2016). De novo discovery of neuropeptides in the genomes of parasitic flatworms using a novel comparative approach. *Int. J. Parasitol.* 46, 709–721.
- Lambert, S.A., Yang, A.W.H., Sasse, A., Cowley, G., Albu, M., Caddick, M.X., Morris, Q.D., Weirauch, M.T., and Hughes, T.R. (2019). Similarity regression predicts evolution of transcription factor sequence specificity. *Nat. Genet.* 51, 981–989.
- Lanjuin, A., VanHoven, M.K., Bargmann, C.I., Thompson, J.K., and Sengupta, P. (2003). *Otx/otd* homeobox genes specify distinct sensory neuron identities in *C. elegans*. *Dev. Cell* 5, 621–633.
- Lesch, B.J., Gehrke, A.R., Bulyk, M.L., and Bargmann, C.I. (2009). Transcriptional regulation and stabilization of left-right neuronal identity in *C. elegans*. *Genes Dev.* 23, 345–358.
- Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25, 1754–1760.
- Li, L., Veksler-Lublin, I., and Zinov'yeva, A. (2019). HRPK-1, a conserved KH-domain protein, modulates microRNA activity during *Caenorhabditis elegans* development. *PLoS Genet.* 15, e1008067.
- Liu, X., Long, F., Peng, H., Aerni, S.J., Jiang, M., Sánchez-Blanco, A., Murray, J.I., Preston, E., Mericle, B., Batzoglou, S., et al. (2009). Analysis of cell fate from single-cell gene expression profiles in *C. elegans*. *Cell* 139, 623–633.
- Lorenzo, R., Onizuka, M., Defrance, M., and Laurent, P. (2020). Combining single-cell RNA-sequencing with a molecular atlas unveils new markers for *Caenorhabditis elegans* neuron classes. *Nucleic Acids Res.* 48, 7119–7134.
- Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* 15, 550.
- Lun, A.T.L., Riesenfeld, S., Andrews, T., Dao, T.P., Gomes, T., and Marioni, J.C.; participants in the 1st Human Cell Atlas Jamboree (2019). EmptyDrops: distinguishing cells from empty droplets in droplet-based single-cell RNA sequencing data. *Genome Biol.* 20, 63.
- Masoudi, N., Tavazoie, S., Glenwinkel, L., Ryu, L., Kim, K., and Hobert, O. (2018). Unconventional function of an *Achaete-Scute* homolog as a terminal selector of nociceptive neuron identity. *PLoS Biol.* 16, e2004979.
- McCarthy, D.J., Campbell, K.R., Lun, A.T.L., and Wills, Q.F. (2017). Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics* 33, 1179–1186.
- McInnes, L., Healy, J., Saul, N., and Großberger, L. (2018). UMAP: Uniform Manifold Approximation and Projection. *arXiv*, arXiv:1802.03426.
- Melentijevic, I., Toth, M.L., Arnold, M.L., Guasp, R.J., Harinath, G., Nguyen, K.C., Taub, D., Parker, J.A., Neri, C., Gabel, C.V., et al. (2017). *C. elegans* neurons jettison protein aggregates and mitochondria under neurotoxic stress. *Nature* 542, 367–371.
- Melkman, T., and Sengupta, P. (2005). Regulation of chemosensory and GABAergic motor neuron development by the *C. elegans* *Aristaless/Arx* homolog *alr-1*. *Development* 132, 1935–1949.
- Milo, R., Kashtan, N., Itzkovitz, S., Newman, M.E.J., and Alon, U. (2003). On the uniform generation of random graphs with prescribed degree sequences. *arXiv*, arXiv:cond-mat/0312028.
- Mirabeau, O., and Joly, J.S. (2013). Molecular evolution of peptidergic signaling systems in bilaterians. *Proc. Natl. Acad. Sci. USA* 110, E2028–E2037.
- Mok, D.Z.L., Sternberg, P.W., and Inoue, T. (2015). Morphologically defined sub-stages of *C. elegans* vulval development in the fourth larval stage. *BMC Dev. Biol.* 15, 26.
- Moresco, J.J., and Koelle, M.R. (2004). Activation of EGL-47, a  $\alpha$ (o)-coupled receptor, inhibits function of hermaphrodite-specific motor neurons to regulate *Caenorhabditis elegans* egg-laying behavior. *J. Neurosci.* 24, 8522–8530.
- Moyle, M.W., Barnes, K.M., Kuchroo, M., Gonopolskiy, A., Duncan, L.H., Sengupta, T., Shao, L., Guo, M., Santella, A., Christensen, R., et al. (2021). Structural and developmental principles of neuropil assembly in *C. elegans*. *Nature* 591, 99–104.
- Murru, L., Moretto, E., Martano, G., and Passafaro, M. (2018). Tetraspanins shape the synapse. *Mol. Cell. Neurosci.* 91, 76–81.
- Nguyen, J.P., Shipley, F.B., Linder, A.N., Plummer, G.S., Liu, M., Setru, S.U., Shaevitz, J.W., and Leifer, A.M. (2016). Whole-brain calcium imaging with cellular resolution in freely behaving *Caenorhabditis elegans*. *Proc. Natl. Acad. Sci. USA* 113, E1074–E1081.
- Norris, A.D., Gao, S., Norris, M.L., Ray, D., Ramani, A.K., Fraser, A.G., Morris, Q., Hughes, T.R., Zhen, M., and Calarco, J.A. (2014). A pair of RNA-binding proteins controls networks of splicing events contributing to specialization of neural cell types. *Mol. Cell* 54, 946–959.
- O'Connor, T.P., Cockburn, K., Wang, W., Tapia, L., Currie, E., and Bamji, S.X. (2009). Semaphorin 5B mediates synapse elimination in hippocampal neurons. *Neural Dev.* 4, 18.
- Ortiz, C.O., Etchberger, J.F., Posy, S.L., Frøkjær-Jensen, C., Lockery, S., Honig, B., and Hobert, O. (2006). Searching for neuronal left/right asymmetry: genome-wide analysis of nematode receptor-type guanylyl cyclases. *Genetics* 173, 131–149.
- Packer, J.S., Zhu, Q., Huynh, C., Sivaramakrishnan, P., Preston, E., Dueck, H., Stefanik, D., Tan, K., Trapnell, C., Kim, J., et al. (2019). A lineage-resolved molecular atlas of *C. elegans* embryogenesis at single-cell resolution. *Science* 365, eaax1971.
- Patrick, R., Humphreys, D., Oshlack, A., Ho, J.W.K., Harvey, R.P., and Lo, K.K. (2019). Sierra: Discovery of differential transcript usage from polyA-captured single-cell RNA-seq data. *BioRxiv*. <https://doi.org/10.1101/13059-020-02071-7>.
- Pavy, N., Rombauts, S., Déhais, P., Mathé, C., Ramana, D.V.V., Leroy, P., and Rouzé, P. (1999). Evaluation of gene prediction software using a genomic data set: application to *Arabidopsis thaliana* sequences. *Bioinformatics* 15, 887–899.
- Pereira, L., Kratsios, P., Serrano-Saiz, E., Sheftel, H., Mayo, A.E., Hall, D.H., White, J.G., LeBoeuf, B., Garcia, L.R., Alon, U., and Hobert, O. (2015). A cellular and regulatory map of the cholinergic nervous system of *C. elegans*. *eLife* 4, e12432.
- Petersen, S.C., Watson, J.D., Richmond, J.E., Sarov, M., Walshall, W.W., and Miller, D.M., 3rd. (2011). A transcriptional program promotes remodeling of GABAergic synapses in *Caenorhabditis elegans*. *J. Neurosci.* 31, 15362–15375.
- Pierce-Shimomura, J.T., Faumont, S., Gaston, M.R., Pearson, B.J., and Lockery, S.R. (2001). The homeobox gene *lim-6* is required for distinct chemosensory representations in *C. elegans*. *Nature* 410, 694–698.
- Pols, M.S., and Klumperman, J. (2009). Trafficking and function of the tetraspanin CD63. *Exp. Cell Res.* 315, 1584–1592.

- Poon, V.Y., Klassen, M.P., and Shen, K. (2008). UNC-6/netrin and its receptor UNC-5 locally exclude presynaptic components from dendrites. *Nature* 455, 669–673.
- Poulin, J.F., Tasic, B., Hjerling-Leffler, J., Trimarchi, J.M., and Awatramani, R. (2016). Disentangling neural cell diversity using single-cell transcriptomics. *Nat. Neurosci.* 19, 1131–1141.
- Qiu, X., Mao, Q., Tang, Y., Wang, L., Chawla, R., Pliner, H.A., and Trapnell, C. (2017a). Reversed graph embedding resolves complex single-cell trajectories. *Nat. Methods* 14, 979–982.
- Qiu, X., Hill, A., Packer, J., Lin, D., Ma, Y.A., and Trapnell, C. (2017b). Single-cell mRNA quantification and differential analysis with Censur. *Nat. Methods* 14, 309–315.
- Raj, B., and Blencowe, B.J. (2015). Alternative Splicing in the Mammalian Nervous System: Recent Insights into Mechanisms and Functional Roles. *Neuron* 87, 14–27.
- Reilly, M.B., Cros, C., Varol, E., Yemini, E., and Hobert, O. (2020). Unique homeobox codes delineate all the neuron classes of *C. elegans*. *Nature* 584, 595–601.
- Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010). edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 26, 139–140.
- Robinson, J.T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E.S., Getz, G., and Mesirov, J.P. (2011). Integrative genomics viewer. *Nat. Biotechnol.* 29, 24–26.
- Sengupta, P., Colbert, H.A., and Bargmann, C.I. (1994). The *C. elegans* gene *odr-7* encodes an olfactory-specific member of the nuclear receptor superfamily. *Cell* 79, 971–980.
- Sengupta, P., Chou, J.H., and Bargmann, C.I. (1996). *odr-10* encodes a seven transmembrane domain olfactory receptor required for responses to the odorant diacetyl. *Cell* 84, 899–909.
- Serrano-Saiz, E., Poole, R.J., Felton, T., Zhang, F., De La Cruz, E.D., and Hobert, O. (2013). Modular control of glutamatergic neuronal identity in *C. elegans* by distinct homeodomain proteins. *Cell* 155, 659–673.
- Shan, G., Kim, K., Li, C., and Walthall, W.W. (2005). Convergent genetic programs regulate similarities and differences between related motor neuron classes in *Caenorhabditis elegans*. *Dev. Biol.* 280, 494–503.
- Shen, K., and Bargmann, C.I. (2003). The immunoglobulin superfamily protein SYG-1 determines the location of specific synapses in *C. elegans*. *Cell* 112, 619–630.
- Siebert, S., Farrell, J.A., Cazet, J.F., Abeykoon, Y., Primack, A.S., Schnitzler, C.E., and Juliano, C.E. (2019). Stem cell differentiation trajectories in *Hydra* resolved at single-cell resolution. *Science* 365, eaav9314.
- Siegenthaler, D., Enneking, E.M., Moreno, E., and Pielage, J. (2015). L1CAM/Neuroglian controls the axon-axon interactions establishing layered and lobular mushroom body architecture. *J. Cell Biol.* 208, 1003–1018.
- Smith, S.J., Sümbül, U., Graybuck, L.T., Collman, F., Seshamani, S., Gala, R., Gliko, O., Elabbady, L., Miller, J.A., Bakken, T.E., et al. (2019). Single-cell transcriptomic evidence for dense intracortical neuropeptide networks. *eLife* 8, e47889.
- Soneson, C., and Delorenzi, M. (2013). A comparison of methods for differential expression analysis of RNA-seq data. *BMC Bioinformatics* 14, 91.
- Spencer, W.C., McWhirter, R., Miller, T., Strasbourger, P., Thompson, O., Hillier, L.W., Waterston, R.H., and Miller, D.M., 3rd. (2014). Isolation of specific neurons from *C. elegans* larvae for gene expression profiling. *PLoS ONE* 9, e112102.
- Sperry, R.W. (1963). Chemoaffinity in the orderly growth of nerve fiber patterns and connections. *Proc. Natl. Acad. Sci. USA* 50, 703–710.
- Stefanakis, N., Carrera, I., and Hobert, O. (2015). Regulatory Logic of Pan-Neuronal Gene Expression in *C. elegans*. *Neuron* 87, 733–750.
- Stuart, T., Butler, A., Hoffman, P., Hafemeister, C., Papalexi, E., Mauck, W.M., 3rd, Hao, Y., Stoeckius, M., Smibert, P., and Satija, R. (2019). Comprehensive Integration of Single-Cell Data. *Cell* 177, 1888–1902.e21.
- Sulston, J.E., and Horvitz, H.R. (1977). Post-embryonic cell lineages of the nematode, *Caenorhabditis elegans*. *Dev. Biol.* 56, 110–156.
- Sulston, J.E., Schierenberg, E., White, J.G., and Thomson, J.N. (1983). The embryonic cell lineage of the nematode *Caenorhabditis elegans*. *Dev. Biol.* 100, 64–119.
- Swoboda, P., Adler, H.T., and Thomas, J.H. (2000). The RFX-type transcription factor DAF-19 regulates sensory neuron cilium formation in *C. elegans*. *Mol. Cell* 5, 411–421.
- Tamburino, A.M., Ryder, S.P., and Walhout, A.J.M. (2013). A compendium of *Caenorhabditis elegans* RNA binding proteins predicts extensive regulation at multiple levels. *G3 (Bethesda)* 3, 297–304.
- Tasic, B., Menon, V., Nguyen, T.N., Kim, T.K., Jarsky, T., Yao, Z., Levi, B., Gray, L.T., Sorensen, S.A., Dolbeare, T., et al. (2016). Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nat. Neurosci.* 19, 335–346.
- Thompson, M., Bixby, R., Dalton, R., Vandenburg, A., Calarco, J.A., and Norris, A.D. (2019). Splicing in a single neuron is coordinately controlled by RNA binding proteins and transcription factors. *eLife* 8, e46726.
- Tomioka, M., Naito, Y., Kuroyanagi, H., and Iino, Y. (2016). Splicing factors control *C. elegans* behavioural learning in a single neuron by producing DAF-2c receptor. *Nat. Commun.* 7, 11645.
- Tourasse, N.J., Millet, J.R.M., and Dupuy, D. (2017). Quantitative RNA-seq meta-analysis of alternative exon usage in *C. elegans*. *Genome Res.* 27, 2120–2128.
- Tran, T.S., Rubio, M.E., Clem, R.L., Johnson, D., Case, L., Tessier-Lavigne, M., Hagan, R.L., Ginty, D.D., and Kolodkin, A.L. (2009). Secreted semaphorins control spine distribution and morphogenesis in the postnatal CNS. *Nature* 462, 1065–1069.
- Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokharel, P., Li, S., Morse, M., Lennon, N.J., Livak, K.J., Mikkelsen, T.S., and Rinn, J.L. (2014). The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol.* 32, 381–386.
- Treisman, J., Gönczy, P., Vashishtha, M., Harris, E., and Desplan, C. (1989). A single amino acid can determine the DNA binding specificity of homeodomain proteins. *Cell* 59, 553–562.
- Troemel, E.R., Sagasti, A., and Bargmann, C.I. (1999). Lateral signaling mediated by axon contact and calcium entry regulates asymmetric odorant receptor expression in *C. elegans*. *Cell* 99, 387–398.
- Tursun, B., Patel, T., Kratsios, P., and Hobert, O. (2011). Direct conversion of *C. elegans* germ cells into specific neuron types. *Science* 331, 304–308.
- van den Brink, S.C., Sage, F., Vértessy, Á., Spanjaard, B., Peterson-Maduro, J., Baron, C.S., Robin, C., and van Oudenaarden, A. (2017). Single-cell sequencing reveals dissociation-induced gene expression in tissue subpopulations. *Nat. Methods* 14, 935–936.
- Varadan, V., Miller, D.M., 3rd, and Anastassiou, D. (2006). Computational inference of the molecular logic for synaptic connectivity in *C. elegans*. *Bioinformatics* 22, e497–e506.
- Venkatachalam, V., Ji, N., Wang, X., Clark, C., Mitchell, J.K., Klein, M., Tabone, C.J., Florman, J., Ji, H., Greenwood, J., et al. (2016). Pan-neuronal imaging in roaming *Caenorhabditis elegans*. *Proc. Natl. Acad. Sci. USA* 113, E1082–E1088.
- Vidal, B., Aghayeva, U., Sun, H., Wang, C., Glenwinkel, L., Bayer, E.A., and Hobert, O. (2018). An atlas of *Caenorhabditis elegans* chemoreceptor expression. *PLoS Biol.* 16, e2004218.
- Von Stetina, S.E., Fox, R.M., Watkins, K.L., Starich, T.A., Shaw, J.E., and Miller, D.M., 3rd. (2007). UNC-4 represses CEH-12/HB9 to specify synaptic inputs to VA motor neurons in *C. elegans*. *Genes Dev.* 21, 332–346.
- Vuong, C.K., Black, D.L., and Zheng, S. (2016). The neurogenetics of alternative splicing. *Nat. Rev. Neurosci.* 17, 265–281.
- Wang, T., Li, B., Nelson, C.E., and Nabavi, S. (2019). Comparative analysis of differential gene expression analysis tools for single-cell RNA sequencing data. *BMC Bioinformatics* 20, 40.

- Weirauch, M.T., Yang, A., Albu, M., Cote, A.G., Montenegro-Montero, A., Drewe, P., Najafabadi, H.S., Lambert, S.A., Mann, I., Cook, K., et al. (2014). Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* 158, 1431–1443.
- White, J.G., Southgate, E., Thomson, J.N., and Brenner, S. (1986). The structure of the nervous system of the nematode *Caenorhabditis elegans*. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 314, 1–340.
- Witvliet, D., Mulcahy, B., Mitchell, J.K., Meirovitch, Y., Berger, D.R., Wu, Y., Liu, Y., Koh, W.X., Parvathala, R., Holmyard, D., et al. (2020). Connectomes across development reveal principles of brain maturation in *C. elegans*. *bioRxiv*. <https://doi.org/10.1101/2020.04.30.066209>.
- Yemini, E., Lin, A., Nejatbakhsh, A., Varol, E., Sun, R., Mena, G.E., Samuel, A.D.T., Paninski, L., Venkatachalam, V., and Hobert, O. (2021). NeuroPAL: A Multicolor Atlas for Whole-Brain Neuronal Identification in *C. elegans*. *Cell* 184, 272–288.e11.
- Young, M.D., and Behjati, S. (2020). SoupX removes ambient RNA contamination from droplet-based single-cell RNA sequencing data. *Gigascience* 9, 1–10.
- Yu, S., Avery, L., Baude, E., and Garbers, D.L. (1997). Guanylyl cyclase expression in specific sensory neurons: a new family of chemosensory receptors. *Proc. Natl. Acad. Sci. USA* 94, 3384–3387.
- Zeisel, A., Muñoz-Manchado, A.B., Codeluppi, S., Lönnerberg, P., La Manno, G., Juréus, A., Marques, S., Munguba, H., He, L., Betsholtz, C., et al. (2015). Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* 347, 1138–1142.
- Zhang, S., Banerjee, D., and Kuhn, J.R. (2011). Isolation and culture of larval cells from *C. elegans*. *PLoS ONE* 6, e19505.
- Zheng, Y., Brockie, P.J., Mellem, J.E., Madsen, D.M., and Maricq, A.V. (1999). Neuronal control of locomotion in *C. elegans* is modified by a dominant mutation in the GLR-1 ionotropic glutamate receptor. *Neuron* 24, 347–361.
- Zhou, X., and Bessereau, J.L. (2019). Molecular Architecture of Genetically-Tractable GABA Synapses in *C. elegans*. *Front. Mol. Neurosci.* 12, 304.
- Zhu, Y., Sousa, A.M.M., Gao, T., Skarica, M., Li, M., Santpere, G., Esteller-Cuatala, P., Juan, D., Ferrández-Peral, L., Gulden, F.O., et al. (2018). Spatiotemporal transcriptomic divergence across human and macaque brain development. *Science* 362, eaat8077.



## STAR★METHODS

### KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
<b>Bacterial and virus strains</b>		
<i>E. coli</i> : OP50	Caenorhabditis Genetics Center	WormBase: OP50; WormBase: WBStrain00041969
<i>E. coli</i> : Na22	Caenorhabditis Genetics Center	WormBase: Na22; WormBase: WBStrain00041948
<b>Chemicals, peptides, and recombinant proteins</b>		
TRIzol LS	Ambion	Cat#10296010
Pronase, Protease from <i>Streptomyces griseus</i>	Sigma-Aldrich	Cat#P8811
<b>Critical commercial assays</b>		
10x Chromium Single Cell 3' GEM, Library & Gel Bead Kit v3	10x Genomics	Cat#1000075
High sensitivity DNA reagents (used with Agilent Bioanalyzer 2100 system)	Agilent Technologies	Cat#5067-4626
Novaseq 6000 S4 150bp PE reads	Illumina	Cat#20012866
Agilent RNA 6000 Pico Reagents (used with Agilent 2100 bioanalyzer system)	Agilent Technologies	Cat# 5067-1513
Phase Lock Gel-Heavy Tubes	Quantabio	Cat#2302830
RNA Clean and Concentrator Kit	Zymo Research	Cat#R1015
<b>Deposited data</b>		
Single cell RNA-Seq data generated in this study	This study	GEO: GSE136049
Bulk RNA-Seq data generated in this study	This study	GEO: GSE169137
<b>Experimental models: organisms/strains</b>		
<i>C. elegans</i> : Strain N2	Caenorhabditis Genetics Center	WormBase: N2; WormBase: WBStrain00000001
OH10689 <i>otIs355 [rab-3(prom1)::2xNLS-TagRFP]</i> IV	<a href="#">Stefanakis et al., 2015</a>	OH10689
EG1285 <i>lin-15B&amp;lin-15A(n765); oxIs12 [unc-47p::GFP + lin-15(+)]</i> X	Caenorhabditis Genetics Center	EG1285
NC3582 <i>oxIs12 [unc-47p::GFP + lin-15(+)]</i> X; <i>otIs355 [rab-3(prom1)::2xNLS-TagRFP]</i> IV	This study	NC3582
VM484 <i>akIs3 [nmr-1p::GFP + lin-15(+)]</i> V	<a href="#">Zheng et al., 1999</a>	VM484
NC3572 <i>akIs3 [nmr-1p::GFP + lin-15(+)]</i> V; <i>otIs355 [rab-3(prom1)::2xNLS-TagRFP]</i> IV	This study	NC3572
OH9625 <i>otIs292 [eat-4::mCherry + rol-6(su1006)]</i>	<a href="#">Tursun et al., 2011</a>	OH9625
OH11746 <i>pha-1(e2123) III; otIs447 [unc-3p::mCherry + pha-1(+)]</i> IV	<a href="#">Kerk et al., 2017</a>	OH11746
OH11157 <i>pha-1(e2123) III; otIs393 [ift-20::NLS-TagRFP + pha-1(+)]</i>	<a href="#">Masoudi et al., 2018</a>	OH11157
OH13470 <i>him-5(e1490) V; otIs354 [cho-1(fosmid)::SL2::YFP::H2B]</i>	<a href="#">Pereira et al., 2015</a>	OH13470
NC3579 <i>otIs354 [cho-1(fosmid)::SL2::YFP::H2B]; otIs355 [rab-3(prom1)::2xNLS-TagRFP]</i> IV	This study	NC3579
NC3580 <i>zdlIs13 [tph-1::GFP]</i> IV; <i>hplIs202 [ceh-10p::GFP + lin-15(+)]</i>	This study	NC3580

(Continued on next page)

**Continued**

REAGENT or RESOURCE	SOURCE	IDENTIFIER
CZ631 <i>juls14</i> [ <i>acr-2::GFP + lin-15(+)</i> ] IV	Hallam et al., 2000	CZ631
RW10754 <i>stls10447</i> [ <i>ceh-34p::HIS-24::mCherry + unc-119(+)</i> ]	Liu et al., 2009	RW10754
NW1229 <i>dpy-20(e1362)</i> IV; <i>evls111</i> [ <i>F25B3.3::GFP + dpy-20(+)</i> ]	Caenorhabditis Genetics Center	NW1229
NC3583 <i>stls10447</i> [ <i>ceh-34p::HIS-24::mCherry + unc-119(+)</i> ]; <i>evls111</i> [ <i>F25B3.3::GFP + dpy-20(+)</i> ]	This study	NC3583
OH15430 <i>pha-1(e2123)</i> III; <i>otls669</i> [NeuroPAL 15] V	Yemini et al., 2021	OH15340
OH16474 <i>pha-1(e2123)</i> III; <i>otEx7567</i> [ <i>nlp-56::GFP + pha-1 (+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	OH16474
OH16475 <i>pha-1(e2123)</i> III; <i>otEx7568</i> [ <i>nlp-17::GFP + pha-1 (+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	OH16475
OH16469 <i>pha-1(e2123)</i> III; <i>otEx7652</i> [ <i>flp-33::GFP + pha-1 (+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	OH16469
OH16630 <i>pha-1(e2123)</i> III; <i>otEx7597</i> [ <i>nlp-42::GFP + pha-1 (+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	OH16630
OH16636 <i>pha-1(e2123)</i> III; <i>otEx7603</i> [ <i>nlp-52::GFP + pha-1 (+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	OH16636
CX5974 <i>kyls262</i> [ <i>unc-86::myr-GFP + odr-1::RFP</i> ] IV	Caenorhabditis Genetics Center	CX5974
NC3636 <i>hdls1</i> [ <i>unc-53p::GFP + rol-6(su1006)</i> ]; <i>otls355</i> [ <i>rab-3prom1::2xNLS-tagRFP</i> ] IV	This study	NC3636
OH16003 <i>otls742</i> [ <i>nlp-13p::GFP + lin-15(+)</i> ]	This study	OH16003
PS3504 <i>unc-119(ed4)</i> ; <i>syls54</i> [ <i>ceh-2::GFP + unc-119(+)</i> ]	Inoue et al., 2002	PS3504
OH16144 <i>nls175</i> [ <i>ceh-28p::4xNLS-GFP + lin-15(+)</i> ]	Hirose et al., 2010	OH16144
NC3635 <i>egls1</i> [ <i>dat-1p::GFP</i> ]; <i>uls152</i> [ <i>mec-3p::RFP</i> ]; <i>kyEx1162</i> [ <i>gcy-35p::GFP</i> ]	This study	NC3635
NC3523 <i>wdls90</i> [ <i>unc-4c::GFP</i> ]	This study	NC3523
NC3685 <i>pha-1(e2123)</i> III; <i>wpEx389</i> [C39H7.2::3xNLS-GFP + <i>pha-1 (+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	NC3685
NC3686 <i>pha-1(e2123)</i> III; <i>wpEx403</i> [ <i>nhr-236::3xNLS-GFP + pha-1 (+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	NC3686
Nspc-1; NeuroPAL	This study	Nspc-1; NeuroPAL
NC3687 <i>wgls707</i> [ <i>sptf-1::TY1::EGFP::3xFLAG + unc-119(+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	NC3687
NC3688 <i>wgls643</i> [ <i>hlh-17::TY1::EGFP::3xFLAG + unc-119(+)</i> ]; <i>otls669</i> [NeuroPAL 15] V	This study	NC3688
RW11595 <i>unc-119(tm4063)</i> ; <i>stls11595</i> [ZK930.3b::H1-wCherry + <i>unc-119(+)</i> ]	Caenorhabditis Genetics Center	RW11595
OH14973 <i>pha-1(e2123)</i> ; <i>otEx6966</i> [ <i>srw-119prom::GFP, pha-1(+)</i> ]	Vidal et al., 2018	OH14973
NC1750 (KM173[ <i>opt-3::GFP[pRF4]</i> ]; <i>hdls32</i> [ <i>glr-1::DsRed2</i> ])	This study	NC1750
ZM9592 <i>hpls670</i> [ <i>pnmr-1::GFP ZF; pglr-5::ZIF-1::SL2::wCherry; lin-15(+)</i> ]	This study	ZM9592
PY10421 [ <i>gpa-4p(d6)::myrGFP2.1</i> ]	This study	PY10421
CX3553 <i>lin-15(n765)</i> ; <i>kyls104</i> [ <i>str-1p::GFP, lin-15(+)</i> ]	This study	CX3553

(Continued on next page)

**Continued**

REAGENT or RESOURCE	SOURCE	IDENTIFIER
NC3182 ( <i>otls138 [ser2prom3::GFP + rol-6]; otls396 [ace-1prom2::NLS::TagRFP]; otls181 [dat-1::mCherry]</i> )	This study	NC3182
NC3296 ( <i>juls223 [pttr-39::mCherry; ptx-3::GFP]; ynl37 [flp-13::GFP]</i> )	This study	NC3296
PHX2805 <i>nlp-51(syb2085[nlp-51::T2A::3xNLS::GFP])</i>	This study	PHX2805
PHX2658 <i>flp-1(syb2658[flp-1::T2A::3xNLS::GFP])</i>	This study	PHX2658
<b>Oligonucleotides</b>		
Ovation SoLo RNA-Seq System with Custom AnyDeplete for the depletion of <i>C. elegans</i> rRNA	Tecan Genomics	Cat#30185717
<b>Software and algorithms</b>		
Cellranger version 3.1.0	10x Genomics	<a href="https://support.10xgenomics.com/single-cell-gene-expression/software/downloads/latest">https://support.10xgenomics.com/single-cell-gene-expression/software/downloads/latest</a>
Background correction and merging dataset code	This study	<a href="https://github.com/cengenproject/Initial_single_cell_analysis">https://github.com/cengenproject/Initial_single_cell_analysis</a>
Thresholding code	This study	<a href="https://github.com/cengenproject/Thresholding_sc">https://github.com/cengenproject/Thresholding_sc</a>
Splicing analysis code	This study	<a href="https://github.com/cengenproject/splicing">https://github.com/cengenproject/splicing</a>
Connectivity analysis code	This study	<a href="https://github.com/cengenproject/connectivity_analysis">https://github.com/cengenproject/connectivity_analysis</a>
CeNGENApp code	This study	<a href="https://github.com/cengenproject/CengenApp">https://github.com/cengenproject/CengenApp</a>
R version 3.6.3	R-CRAN	<a href="https://www.r-project.org">https://www.r-project.org</a>
R Studio version 1.2.1335	RStudio	<a href="https://www.rstudio.com">https://www.rstudio.com</a>
MATLAB	MathWorks	R2019b
R package DropletUtils version 1.6.1	R Bioconductor; (Lun et al., 2019)	<a href="https://www.bioconductor.org/packages/release/bioc/html/DropletUtils.html">https://www.bioconductor.org/packages/release/bioc/html/DropletUtils.html</a>
R package Seurat version 3.1.5	Github; (Stuart et al., 2019)	<a href="https://github.com/satijalab/seurat">https://github.com/satijalab/seurat</a>
R package scater version 1.14.6	R Bioconductor; (McCarthy et al., 2017)	<a href="https://bioconductor.org/packages/release/bioc/html/scater.html">https://bioconductor.org/packages/release/bioc/html/scater.html</a>
R package monocle3 version 0.2.2	Github; (Qiu et al., 2017a)	<a href="https://cole-trapnell-lab.github.io/monocle3/docs/installation/">https://cole-trapnell-lab.github.io/monocle3/docs/installation/</a>
R package igraph version 1.2.5	R CRAN	<a href="https://cran.r-project.org/web/packages/igraph/index.html">https://cran.r-project.org/web/packages/igraph/index.html</a>
R package ggpubr version 0.4.0	R CRAN	<a href="https://cran.r-project.org/web/packages/ggpubr/index.html">https://cran.r-project.org/web/packages/ggpubr/index.html</a>
R package pheatmap version 1.0.12	R CRAN	<a href="https://cran.r-project.org/web/packages/pheatmap/index.html">https://cran.r-project.org/web/packages/pheatmap/index.html</a>
R package ggplot2 version 3.3.2	R CRAN	<a href="https://cran.r-project.org/web/packages/ggplot2/index.html">https://cran.r-project.org/web/packages/ggplot2/index.html</a>
R package SoupX version 1.4.5	Github; Young and Behjati, 2020	<a href="https://github.com/constantAmateur/SoupX">https://github.com/constantAmateur/SoupX</a>
R package boot version 1.3-24	R-CRAN; Canty and Ripley, 2019	<a href="https://cran.r-project.org/web/packages/boot/">https://cran.r-project.org/web/packages/boot/</a>
SAMtools version 1.9	Github	<a href="https://github.com/samtools/">https://github.com/samtools/</a>
STAR version 2.7.0	Github	<a href="https://github.com/alexdobin/STAR">https://github.com/alexdobin/STAR</a>
SubRead version 1.6.4	SourceForge	<a href="http://subread.sourceforge.net/">http://subread.sourceforge.net/</a>
R package edgeR version 3.28.1	R Bioconductor	<a href="https://bioconductor.org/packages/release/bioc/html/edgeR.html">https://bioconductor.org/packages/release/bioc/html/edgeR.html</a>
SplAdder	Github; Kahles et al., 2016	<a href="https://github.com/ratschlab/spladder">https://github.com/ratschlab/spladder</a>
FIRE	Elemento et al., 2007	<a href="https://tavazoilab.c2b2.columbia.edu/FIRE/">https://tavazoilab.c2b2.columbia.edu/FIRE/</a>

(Continued on next page)

## Continued

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Other		
Confocal Laser Scanning Microscope	Zeiss	LSM 880
Confocal Laser Scanning Microscope	Nikon	A1R
CeNGEN website	This study	<a href="https://www.cengen.org">https://www.cengen.org</a>
CeNGENApp web application	This study	<a href="https://cengen.shinyapps.io/CengenApp/">https://cengen.shinyapps.io/CengenApp/</a>

## RESOURCE AVAILABILITY

### Lead contact

Requests for resources and reagents should be directed to the lead contact, David Miller ([david.miller@vanderbilt.edu](mailto:david.miller@vanderbilt.edu))

### Materials availability

The strains generated in this study are available at the *Caenorhabditis* Genetics Center or by request from the lead contact.

### Data and code availability

The raw data are available at GEO (single cell data, GEO: GSE136049; bulk sequence data, GEO: GSE169137). The full and neuron only datasets are available at <https://www.cengen.org>. Analysis code is available at github <https://github.com/cengenproject>.

## EXPERIMENTAL MODEL AND SUBJECT DETAILS

### Preparation of larvae and dissociation

Worms were grown on 8P nutrient agar 150 mm plates seeded with *E. coli* strain NA22. To obtain synchronized cultures of L4 worms, embryos obtained by hypochlorite treatment of adult hermaphrodites were allowed to hatch in M9 buffer overnight (16–23 hours at 20°C) and then grown on NA22-seeded plates for 45–48 hours at 23°C. The developmental age of each culture was determined by scoring vulval morphology (> 75 worms) (Mok et al., 2015). Single cell suspensions were obtained as described (Kaletsky et al., 2016; Spencer et al., 2014; Zhang et al., 2011) with some modifications. Worms were collected and separated from bacteria by washing twice with ice-cold M9 and centrifuging at 150 *rcf*. for 2.5 minutes. Worms were transferred to a 1.6 mL centrifuge tube and pelleted at 16,000 *rcf*. for 1 minute. 250  $\mu$ L pellets of packed worms were treated with 500  $\mu$ L of SDS-DTT solution (20 mM HEPES, 0.25% SDS, 200 mM DTT, 3% sucrose, pH 8.0) for 2–4 minutes. In initial experiments, we noted that SDS-DTT treatment for 2 minutes was sufficient to dissociate neurons from the head and tail, but longer times were required for effective dissociation of neurons in the mid-body and ventral nerve cord. The duration of SDS-DTT was therefore selected based on the cells targeted in each experiment. For example, NC3582, OH11746, and *juls14* L4 larvae were treated for 4 minutes to ensure dissociation and release of ventral cord motor neurons. NC3579, NC3580 and NC3636 L4 larvae were treated with SDS-DTT for 3 minutes. All other strains were incubated in SDS-DTT for 2 minutes. Following SDS-DTT treatment, worms were washed five times by diluting with 1 mL egg buffer and pelleting at 16,000 *rcf*. for 30 s. Worms were then incubated in pronase (15 mg/mL, Sigma-Aldrich P8811, diluted in egg buffer) for 23 minutes. During the pronase incubation, the solution was triturated by pipetting through a P1000 pipette tip for four sets of 80 repetitions. The status of dissociation was monitored under a fluorescence dissecting microscope at 5-minute intervals. The pronase digestion was stopped by adding 750  $\mu$ L L-15 media supplemented with 10% fetal bovine serum (L-15-10), and cells were pelleted by centrifuging at 530 *rcf*. for 5 minutes at 4 C. The pellet was resuspended in L-15-10, and single-cells were separated from whole worms and debris by centrifuging at 100 *rcf*. for 2 minutes at 4 C. The supernatant was then passed through a 35-micron filter into the collection tube. The pellet was resuspended a second time in L-15-10, spun at 100 *rcf*. for 2 minutes at 4 C, and the resulting supernatant was added to the collection tube.

## METHOD DETAILS

### FACS isolation of neuron types for RNA-Seq

Fluorescence Activated Cell Sorting (FACS) was performed on a BD FACSAria III equipped with a 70-micron diameter nozzle. DAPI was added to the sample (final concentration of 1  $\mu$ g/mL) to label dead and dying cells. To prepare samples for scRNA-sequencing, our general strategy used fluorescent reporter strains to isolate subgroups of cells. For example, we used an *eat-4::mCherry* reporter (OH9625) to target glutamatergic neurons and an *ift-20::NLS-TagRFP* reporter (OH11157) to label ciliated sensory neurons. We used an intersectional labeling strategy with a nuclear-localized pan-neural marker (*otIs355 [rab-3(prom1)::2xNLS-TagRFP]* IV) to exclude cell fragments labeled with cytosolic GFP markers (NC3582). In other cases, we used an intersectional strategy to exclude non-neuronal cells. For example, *stIs10447 [ceh-34p::HIS-24::mCherry]* is expressed in pharyngeal muscles, pharyngeal neurons and



coelomocytes. To target pharyngeal neurons, we generated strain NC3583 by crossing *stIs10447 [ceh-34p::HIS-24::mCherry]* with the pan-neural GFP marker *evIs111* to isolate cells that were positive for both mCherry and GFP. Non-fluorescent N2 (wild-type reference strain) (Brenner, 1974) standards and single-color controls (in the case of intersectional labeling approaches) were used to set gates to exclude auto-fluorescent cells and to compensate for bleed-through between fluorescent channels. For two experiments, single-cell suspensions from separate strains were combined (OH16003 plus PS3504 and *nIs175*, NC3635 plus NC3532) prior to FACS. In some cases, we expanded FACS gates to encompass a wide range of fluorescent intensities to ensure capture of targeted cell types. This less stringent approach may contribute to the presence of non-neuronal cells in our dataset (see Results and Discussion). Cells were sorted under the “4-way Purity” mask.

For 10X Genomics single-cell experiments, sorted cells were collected into L-15-33 (L-15 medium containing 33% fetal bovine serum), concentrated by centrifugation at 500 *rcf* for 12 minutes at 4°C, and counted on a hemocytometer. Single-cell suspensions used for 10x Genomics single-cell sequencing ranged from 300–900 cells/ $\mu$ L.

For bulk RNA-sequencing of individual cell types, sorted cells were collected directly into TRIzol LS. At ~15-minute intervals during the sort, the sort was paused, and the collection tube with TRIzol was inverted 3–4 times to ensure mixing. Cells in TRIzol LS were stored at –80°C for RNA extractions (see below).

### Single-cell RNA sequencing

Each sample (targeting 5,000 or 10,000 cells per sample) was processed for single cell 3' RNA sequencing utilizing the 10X Chromium system. Libraries were prepared using P/N 1000075, 1000073, and 120262 following the manufacturer's protocol. The libraries were sequenced using the Illumina NovaSeq 6000 with 150 bp paired end reads. Real-Time Analysis software (RTA, version 2.4.11; Illumina) was used for base calling and analysis was completed using 10X Genomics Cell Ranger software (v3.1.0). Most samples were processed with 10x Genomics v2 Chemistry, except for samples from *juls14*, NC3583, NC3636, CX5974, OH16003, PS3504, *nIs175*, NC3635 and NC3532, which were processed with v3 Chemistry. Detailed experimental information is found in Table S1.

### Single-cell RNA-Seq Mapping

Reads were mapped to the *C. elegans* reference transcriptome from WormBase, version WS273. Due to the possibility that 3' untranslated region (UTR) annotations in the reference transcriptome may be too short (Packer et al., 2019), we dynamically extended the 3' UTR of each gene to its optimal length, thereby enabling the additional mapping of reads to the 3' extremity of the gene body. We generated eight versions of gene annotations based on WormBase WS273 annotation, with 3' UTRs in each version elongated by 50, 100, 150, 200, 250, 300, 400 and 500 base pairs (bps), respectively. Elongation of genes which overlapped with other genes during the extension process was terminated before encountering an adjacent exon. Subsequently, eight custom genome indexes, which respectively combined the *C. elegans* WS273 reference genome with the eight extended gene annotation versions, were generated using CellRanger (version 3.1.0).

All sequenced reads from each of the 17 single-cell samples were mapped to the eight reference genomes using the CellRanger pipeline. We next selected the best UTR extension length of each annotated gene independently for the 17 samples, as a number of genes were heavily enriched in specific samples. First, we calculated the total number of mapped reads for each of the expressed genes in each sample, resulting in eight mapped-read values representing the eight gene annotation versions. To discard the UTR extension intervals which harbor sparse additional reads, as well as to allow for the intervals which harbor fewer reads but are surrounded by read-enriched intervals, we took advantage of the trimming algorithm in Burrows-Wheeler Alignment (Li and Durbin, 2009) to find the best extension. Specifically, a cutoff of 20 reads was applied to each extension interval (50, 50, 50, 50, 50, 100, and 100 bps). Cumulative sums from 3' to 5' end were then calculated after subtracting the cutoff in each interval, and the smallest sum of less than 0 was located as the trimming point for a given sample. Considering all 17 samples, the trimming point agreed by most samples (or at least two samples if one gene is expressed in limited samples) was chosen as the ultimate one. Consequently, we extended the UTRs for 1,012 *C. elegans* genes, encompassing 40, 216, 175, 113 and 468 genes with UTRs extended by 150, 200, 250, 300 and 400 bps at the 3' end, respectively. Lastly, with the gene annotation file containing the optimal extension length for each gene, we remapped and quantified the gene expression in all 17 samples using CellRanger. The gene annotation file is available at GEO: GSE136049.

### Downstream Processing

We distinguished cells from empty droplets, corrected background RNA expression and generated quality control metrics for each sample independently, then merged the files together into one dataset. The default barcode filtering algorithm in CellRanger can fail to capture cells in some conditions, especially with cells with variable sizes and RNA content (Lun et al., 2019). Neurons in particular tend to have lower UMI counts than other cell types and can be missed by the default algorithm (Packer et al., 2019). We therefore used the EmptyDrops method (with a threshold of 50 UMIs for determining empty droplets) from the R package DropletUtils (Lun et al., 2019) to determine which droplets contained cells. This approach detected significantly more cells than the CellRanger method, and we were able to confidently annotate these additional cells as neurons.

The SoupX R package (Young and Behjati, 2020) was used to correct for background RNA. We used a more conservative threshold for determining background RNA for SoupX than for EmptyDrops to exclude low-quality cells in the background correction. We therefore set a threshold of droplets with fewer than 25 UMIs to estimate the background RNA. Genes with patterns of strong expression in restricted sets of cells (from the literature or from preliminary clustering analysis for each single-cell experiment) were selected for

each dataset (Table S1). SoupX uses these genes, preliminary clustering, and the calculated background RNA profile (from droplets with fewer than 25 UMIs) to estimate the percent of contamination in each sample. The estimated background contamination ranged from 4.15%–13.56%, with a mean of 8.01%. For the *ceh-28\_dat-1* experiment, no combination of genes tested resulted in satisfactory performance, so the contamination was set manually to 10.00%. SoupX uses the calculated contamination level to correct the expression of genes that are abundant in the background RNA profile, and returns a corrected gene by cell count matrix. The background corrected count matrices produced by SoupX were rounded to integer counts and used for subsequent downstream processing.

Following background correction, quality control metrics were calculated for each dataset with the R package *scater* (McCarthy et al., 2017), using the percentage of UMIs from the mitochondrial genes *nduo-1*, *nduo-2*, *nduo-3*, *nduo-4*, *nduo-5*, *nduo-6*, *ctc-1*, *ctc-2*, *ctc-3*, *ndfl-4*, *atp-6*, and *ctb-1*. Droplets with greater than twenty percent of UMIs coming from mitochondrial genes were removed. Datasets from individual experiments were merged using Seurat (v3) (Stuart et al., 2019). Genes detected in fewer than five cells were removed. Log-normalized expression matrices were then used for downstream analysis using monocle (2.99.3), monocle3 (0.2.1) (Cao et al., 2019; Qiu et al., 2017a, 2017b; Trapnell et al., 2014) and Seurat (v3) packages.

### Dimensionality reduction and batch correction

We imported the merged dataset into monocle3, and reduced the dimensionality of the dataset with PCA (135 principal components, based on examination of an elbow plot showing the variance explained by each principal component), followed by the Uniform Manifold Approximation and Projection (UMAP) (Becht et al., 2018; McInnes et al., 2018) algorithm in monocle3 (`reduce_dimension` function, parameters were default other than: `umap.min_dist = 0.3`, `umap.n_neighbors = 75`). We then clustered cells using the leiden algorithm in monocle3 (`res = 3e-4`). Batch correction between experiments was performed using the `align_cds` function (Cao et al., 2019; Haghverdi et al., 2018). We processed the neuron-only dataset with the following parameters (125 PCs, `umap.min_dist = 0.3`, `umap.n_neighbors = 75`, `alignment_k` (for `align_cds`) = 5, clustering resolution 3e-3).

### Cell Identification

We assigned tissue and cell identity to the majority of cells in our dataset based on a manually compiled list of reported gene expression profiles with an average of > 20 molecular markers per neuron type (Hobert et al., 2016), and a recently described protein expression atlas of > 100 homeodomain proteins (Reilly et al., 2020) (Table S1). Most of the neuronal UMAP clusters could be readily assigned to an individual neuron type on the basis of these known markers. We manually excluded clusters we identified as doublets due to co-expression of cell-type specific markers. We manually merged multiple clusters that corresponded to the same neuron type. We noted that coelomocytes were most abundant in experiments using strains expressing mCherry (*otIs292* and *otIs447*). This effect likely results from neurons shedding mCherry+ exophers, which are then taken up by coelomocytes (Melentijevic et al., 2017), causing them to be isolated along with mCherry-labeled neurons.

Some clusters in the initial global dataset appeared to contain multiple closely related neuron types (i.e., cholinergic motor neurons, dopaminergic neurons, oxygen sensing neurons AQR, PQR, URX and pharyngeal neurons). Additional analysis of these separate clusters (i.e., reapplication of PCA, UMAP, and clustering to just these clusters) separated these cell types into individual clusters (Figures 1E and 1F). Finally, we identified separate clusters for the neuron classes RIV and SMD. In both of these instances, however, one of the putative clusters showed strong expression of stress-related transcripts rather than subtype specific markers and therefore likely correspond to a subset of RIV and SMD neurons damaged by the isolation protocol. These two aberrant clusters were excluded from further analyses.

In the complete dataset, cells had a median of 928 UMIs/cell and 328 genes/cell. In the neuron only dataset, neurons had a median of 1033 UMIs/cell and 363 genes/cell. We note that these metrics are lower than generally observed for *Drosophila* or mouse 10X experiments (10X Genomics, 2017; Davie et al., 2018). We believe that this is likely due to the lower RNA content in *C. elegans* neurons (~2  $\mu$ m in diameter) compared to *Drosophila* (2–6  $\mu$ m) or mouse (10–30  $\mu$ m) neurons.

### Neuron network analysis

The neuron network containing all neuron types was constructed on the basis of the transcriptome similarity between each pair of neuron types. We obtained the transcriptional profile of each neuron type by averaging gene expression across all cells within the given type, resulting in the gene expression trajectory for each neuron type. We next calculated transcriptome similarity (after log transformation) as the Pearson correlation coefficient between pairwise neuron types, using 7,390 highly variable genes identified by Seurat based on their variance and mean expression. The neuron network in a graphopt layout was constructed by the package “igraph” (Csardi and Nepusz, 2006) in R using the force-directed graphopt algorithm based on the above similarity matrix.

### Gene expression analyses

Averaged gene expression profiles for each neuron class were generated as described (Cao et al., 2017). Quantitative expression data for a subset of genes are distorted by overexpression from fosmid reporters or co-selectable markers (*lin-15A*, *lin-15B*, *pha-1*, *rol-6*, *unc-119*, *dpy-20*, *cho-1*), the promoter regions used for marking cell types (*unc-53*, *unc-47*, *gcy-35*, *C30A5.16*, *saeg-2*, *F38B6.2*, *C30F8.3*, *cex-1*) or from a gene-specific 3' UTR included in fluorescent reporter constructs (*eat-4*, *unc-54*). These genes are annotated in the CengenApp web application.

For visualization of gene expression data in the web application and for differential gene expression tests, data were imported into Seurat (v3) and raw counts were normalized using the variance stabilizing transformation (VST) implemented in the function `sctransform` with default parameters and regressing out the percent of mitochondrial reads (Hafemeister and Satija, 2019; Stuart et al., 2019). Differential gene expression tests used the Seurat v3 default Wilcoxon rank sum test with default parameters (a gene must be detected in > 10% of the cells in the higher-expressing cluster and have an adjusted p value < 0.05).

### Stress-induced genes

The dissociation procedure used to isolate single cells can induce cellular stress responsive pathways (van den Brink et al., 2017; Kaletsky et al., 2016). To identify likely stress-induced genes, we examined the distribution in our data of a list of 199 stress-induced genes, including heat shock protein (*hsp*) family genes and additional genes from the literature (van den Brink et al., 2017; Brunquell et al., 2016; Kaletsky et al., 2016) (Table S1). 20 of these genes showed abundant and broad expression across the entire nervous system. We generated a stress index for each single cell by calculating the percent of UMIs mapping to these 20 genes. We then tested the correlation of each gene's expression pattern with the stress index to identify additional putative stress-responsive genes. We identified a total of 49 genes featuring correlations > 0.1 with the stress index and which were detected in at least 75 neuron types as likely stress responsive genes (Table S1).

### Thresholding

The wealth of known gene expression data in *C. elegans* from fluorescent reporter strains provides a unprecedented opportunity to set empirical thresholds for our scRNA-Seq data based on ground truth. We first compiled a ground truth dataset of 160 genes with expression patterns across the nervous system previously determined with high confidence fosmid fluorescent reporters, CRISPR strains or other methods (Bhattacharya et al., 2019; Harris et al., 2020; Reilly et al., 2020; Stefanakis et al., 2015; Yemini et al., 2021) (Table S2). For each gene, we then aggregated expression across the single cells corresponding to each neuron type and calculated several metrics, including the total UMI count, the number of single cells of each neuron type in which each gene was detected with at least one UMI, the proportion of single cells of each neuron type in which gene was detected with at least one UMI and a normalized transcripts per million (TPM) expression value (Packer et al., 2019). We generated receiver operating characteristic (ROC) and precision recall (PR) curves for each metric by thresholding the data across a range of values, and calculated true positive, false positive, and false discovery rates by comparing the single-cell data to the ground truth. We used the area under the curve (AUC) to decide which metric to use for thresholding. The proportion of cells in which a gene was detected performed the best (had the highest AUC) and was thus used to establish gene-level thresholds.

We first set initial thresholds to retain ubiquitously-expressed genes and to remove non-neuronal genes. Genes detected in  $\geq 1\%$  of the cells in every neuron cluster were considered expressed in all neuron types (193 genes), whereas transcripts detected in  $\leq 2\%$  of the cells in every neuron cluster were considered non-neuronal (4806 genes; no genes were detected in  $\geq 1\%$  and  $\leq 2\%$  of the cells in every neuron). As most genes displayed different levels of expression, we found that a single threshold failed to reliably capture expression for all genes. Thus, we applied percentile thresholding for each gene individually. For example, the AFD cluster showed the highest proportion of cells (76.3%, Figure S5A) expressing the homeodomain transcription factor *ttx-1*. For *unc-25*/GAD, the VD\_DD cluster had the highest proportion of cells (94.4%, Figure S5G), whereas for the homeodomain transcription factor *ceh-13*, the DA neuron cluster had the highest proportion (13.4%, not shown). Thresholds were calculated as a fraction of the highest proportion of cells for each individual gene. For example, a threshold of 0.04 results in different absolute cut-offs for each gene. For *ttx-1*, with a highest proportion of 76.3%, we scored *ttx-1* as “not expressed” in clusters in which it was detected in < 3.05% of cells ( $0.04 \times 76.3 = 3.05\%$ ). For *unc-25*, with a highest proportion of expressing cells of 94.4%, we scored *unc-25* as “not expressed” in clusters in which it was detected in < 3.77% of cells ( $0.04 \times 94.4 = 3.77\%$ ). Similarly, and we scored *ceh-13* as “not expressed” in clusters in which it was detected in < 0.536% of cells ( $0.04 \times 13.4 = 0.536\%$ ).

For each threshold percentile, we generated 5,000 stratified bootstraps of the ground truth genes using the R package `boot` (Canty and Ripley, 2019; Davison and Hinkley, 1997) and computed the True Positive Rate (TPR), False Positive Rate (FPR) and False Discovery Rate (FDR) for the entire dataset as well as for each neuron type. We estimated 95% confidence intervals with the adjusted percentile (BCa) method, and plotted the ROC and PR curves (Figures S5C and S5D). Finally, we selected 4 thresholds of increased stringency (1-4, see Table S2 for statistics for each neuron type). Threshold 2 was used for analyses profiling gene expression across all neuron types and across gene families.

### Estimating coverage for individual neurons

We used threshold 2 to model the relationship between the number of cells in each neuron type cluster and the number of genes detected with the expression:

$$G_N = G_{max} * \frac{N_C}{b + N_C} \quad (\text{Eq. 1})$$

Where  $G_N$  is the number of genes detected,  $G_{max}$  is the maximal number of genes detected with an infinite number of cells,  $N_C$  is the number of cells of a given type, and  $b$  is the number of cells at which  $G_N$  = half of  $G_{max}$ . Using 1000 bootstrapped samples, we estimate

6550  $\pm$  7 genes for  $G_{max}$  and 34.22  $\pm$  0.3 for  $b$  (Figure S5I). In other words, this finding suggests that single cell sequencing would detect an average of  $\sim$ 6,500 transcripts per neuron type if an infinite number of cells were sampled and that sampling of  $\sim$ 30 cells/neuron type is sufficient to capture 50% of these genes.

To address the possibility that transcript complexity could vary across neuron types, we used a down-sampling strategy to model the relationship between genes detected versus the number of cells sampled for each neuron class. We performed 60–100 iterations of down-sampling for each neuron type to generate plots of numbers of cells versus numbers of genes for each cell type at threshold 2 (Figure S5K). Fitting Equation 1 to each plot predicts a maximal number of genes detected at an infinite number of cells for each neuron type (Figure S5L; Table S2). Estimates for some neuron types are less confident due to under-sampling of cells. However, we also see a wide range of predicted values among well-represented cell types, suggesting that these estimates could be indicative of biological variation in the genetic complexity of individual neuron types across the nervous system (Table S2).

### Determining distinct combinations of gene sets

Expression matrices of selected gene families from threshold 2 were binarized. Genes were clustered following default parameters in the R package hclust. We determined if neurons expressed a distinct combinatorial code for given gene families by determining whether any two columns (neurons) of the binarized expression matrix were identical. For analyzing expression of gene regulatory families, we treated C2H2 zinc finger proteins as transcription factors and removed them from the list of RNA-binding proteins. We also removed ribosomal proteins from the RNA-binding protein list.

### Connectivity Analysis

To determine neurons postsynaptic to either ACh or glutamate-releasing neurons, we used the *C. elegans* hermaphrodite chemical connectome data from (Cook et al., 2019). For this analysis, we scored synapses as connections detected in more than 3 electron micrograph sections.

### Reporter strains

GFP reporters for the neuropeptide genes *flp-33*, *nlp-17*, *nlp-42*, *nlp-52* and *nlp-56* were created by PCR Fusion (Hobert, 2002) whereby the 5' intergenic region of the gene of interest and the coding sequence of GFP with 3' UTR of *unc-54* were fused in subsequent PCR reactions. We used the entire intergenic region of the genes of interest: 1519 bp for *flp-33* (forward primer: aggaagtgtat aaacttgctgttttaaatg, reverse primer: ggtagggggaccctggaag), 372 bp for *nlp-17* (forward primer: tcactctaaaatatatttcaaaacgattttctgtgc, reverse primer: attttctgtgaaaaagcctgacttttc), 3250 bp for *nlp-42* (forward primer: ttgtctgaaaatatgggtttgcatgg, reverse primer: ttactctgaaaatttgaattttcagattttac), 3731 bp for *nlp-52* (forward primer: ttgcttgcaattttctgaaataagatgg, reverse primer: ttttgggaagaggt acctggaac), and 2954 bp for *nlp-56* (forward primer: gggtcactggaataaatatgcactgtatc, reverse primer: ctggaagaggtgaatcatatggttta-gaag). Reporters were injected directly into NeuroPAL *pha-1* strain (OH15430 *pha-1(e2123)*; *otIs669*[NeuroPAL 15]) (Yemini et al., 2021) as a complex array with OP50 DNA (linearized with *ScaI*) and *pBX [pha-1 (+)]* (Granato et al., 1994) as a co-injection marker. For *flp-33* and *nlp-52*, the reporter, *pBX [pha-1 (+)]* and OP50 DNA were injected at concentrations of 7.75 ng/ $\mu$ l, 6.2 ng/ $\mu$ l, 99.96 ng/ $\mu$ l, respectively. For *nlp-42*, the reporter, *pBX [pha-1 (+)]* and OP50 DNA were injected at 11.80 ng/ $\mu$ l, 8.7 ng/ $\mu$ l and 88.86 ng/ $\mu$ l. For *nlp-17*, the reporter, *pBX [pha-1 (+)]* and OP50 DNA were injected at 10 ng/ $\mu$ l, 6.2 ng/ $\mu$ l and 99.96 ng/ $\mu$ l. For *nlp-56*, the reporter, *pBX [pha-1 (+)]* and OP50 DNA were injected at concentrations of 9.5 ng/ $\mu$ l, 5.2 ng/ $\mu$ l and 94.9 ng/ $\mu$ l. After injection, animals were kept at 25°C for selection of the array positive worms and maintained for at least three generations before imaging (see below). CRISPR reporter strains for *flp-1* and *nlp-51* were generated by engineering a T2A::3xNLS::GFP cassette into the respective gene loci just before the stop codons. The *npsc-1* promoter fusion reporter was constructed using the entire 713 bp intergenic region upstream of *npsc-1* fused driving GFP.

Sequences of C39H7.2 and *nhr-236* were acquired from *C. elegans* BioProject PRJNA13758 browser (via WormBase). We combined 1447 bp upstream of the C39H7.2 sequence (forward primer: Gtatgtgtcgcaggatgac, reverse primer: Gcccatggaagtgtcgaatt) with 2044 bp of *UberPN::3xNLS-intronGFP* (forward primer: CCCAAAGgtatgtttcgaat, reverse primer: AACTGTTTCCTACTAGTCGG) via overlap PCR. For *nhr-236*, we combined 802 bp immediately upstream of the ATG sequence of the first exon of *nhr-236* (forward primer: Tcttgaaggcagcccgatt, reverse primer: Gctctgtgtcggattccgg) with 2044 bp of *UberPN::3xNLS-intronGFP* (primers as above) via overlap PCR. The resulting overlap PCR products were injected with 50 ng/ $\mu$ l of *pha-1* rescue construct *pBX [pha-1 (+)]* and 1Kb+ladder (Promega Corporation, G5711) into GE24 *[pha-1(e2123) III]*. The injected lines were grown at 25 C for selection of the *pha-1+* worms and were maintained for at least five generations before imaging with a Spinning Disk Confocal microscope (Nikon). The images were analyzed using Volocity Imaging Software and also crossed into the NeuroPAL strain *otIs669* to identify the neurons expressing the reporters.

### Imaging

Confocal images were obtained on either a Nikon A1R confocal laser scanning microscope or a Zeiss LSM 880 microscope using 20x or 40x oil immersion objectives. Brightness and contrast adjustments were performed with FIJI.



### RNA Extraction

Cell suspensions in TRIzol LS (stored at  $-80^{\circ}\text{C}$ ) were thawed at room temperature. Chloroform extraction was performed using Phase Lock Gel-Heavy tubes (Quantabio) according to the manufacturer's protocol. The aqueous layer from the chloroform extraction was combined with an equal volume of 100% ethanol and transferred to a Zymo-Spin IC column (Zymo Research). Columns were centrifuged for 30 s at 16,000 *rcf.*, washed with 400  $\mu\text{L}$  of Zymo RNA Prep Buffer and centrifuged for 16,000 *rcf.* for 30 s. Columns were washed twice with Zymo RNA Wash Buffer (700  $\mu\text{L}$ , centrifuged for 30 s, followed by 400  $\mu\text{L}$ , centrifuged for 2 minutes). RNA was eluted by adding 15  $\mu\text{L}$  of DNase/RNase-Free water to the column filter and centrifuging for 30 s. A 2  $\mu\text{L}$  aliquot was submitted for analysis using the Agilent 2100 Bioanalyzer Picochip to estimate yield and RNA integrity and the remainder stored at  $-80^{\circ}\text{C}$ .

### Bulk sequencing and mapping

Each bulk RNA sample was processed for sequencing using the SoLo Ovation Ultra-Low Input RNaseq kit from Tecan Genomics according to manufacturer instruction, modified to optimize rRNA depletion for *C. elegans* (Barrett et al., 2021). Libraries were sequenced on the Illumina HiSeq 2500 with 75 bp paired end reads. Reads were mapped to the *C. elegans* reference transcriptome from WormBase (version WS274) using STAR version 2.7.0. Duplicate reads were removed using SAMtools (version 1.9), and a counts matrix was generated using the featureCounts tool of SubRead (version 1.6.4).

### Comparing scRNA-Seq and bulk RNA data

Differential gene expression comparing sorted cell samples with sorted pan-neuronal samples was performed using TMM-normalized counts in edgeR (version 3.28.1). Two to five replicates per cell type were used in each sample (ASG: 4, AVE: 3, AVG: 3, AWA: 4, AWB: 5, DD: 3, PVD: 2, VD: 4, pan-neuronal: 5). Marker genes from the single cell dataset were selected using a Wilcoxon test in Seurat v3, calling enriched genes by comparing individual neuronal clusters to all other neuronal clusters. Marker genes were defined as genes with a log fold change  $> 2$ , and adjusted p value  $< 0.001$ . To examine marker gene enrichment in each bulk cell type, pairwise Wilcoxon tests were performed in R comparing the corresponding bulk cell type's enrichment against the enrichment in all other bulk cell types.

To compare the overlap of gene detection between bulk and single cell datasets, bulk TMM counts were normalized to gene length, and the true positive rate (TPR) for detecting ground truth markers (see Thresholding) was calculated for a range of length normalized TMM values. At each expression threshold, if  $> 65\%$  of samples showed expression equal to or higher than the threshold, the gene was called expressed. TPR, FPR, and FDR rates were calculated with 5,000 stratified bootstraps of the ground truth genes, which were generated using the R package boot (Canty and Ripley, 2019; Davison and Hinkley, 1997). We used a threshold of 5.7 length normalized TMM, to match the TPR (0.81) of the single cell Threshold 2. To calculate the relationship between single cell cluster size and the overlap between bulk and single cell gene expression, only protein coding genes were considered. Classifications from WormBase were used to define each gene's RNA class.

### Alternative splicing

Alternative splicing events were detected using the software SplAdder (Kahles et al., 2016). The common splicing graph was built based on all 32 individual samples and each pair of neurons was tested for differential use of AS events (with confidence level of 3 and parameters=ignore-mismatches,-validate-sg and sg\_min\_edge\_count = 3). The resulting tables were loaded in R to adjust the p value for multiple testing, and events with FDR  $> 0.1$  were discarded. Sashimi plots for the genes *mca-3* and *mbk-2* were generated using the Integrated Genomics Viewer (Robinson et al., 2011).

For the previously unannotated exons, the splicing graph generated by SplAdder was recovered. It consisted of 197,576 exons; of these, 3,860 were not annotated in WormBase WS274. To avoid counting exons resulting from intron retention events or imprecise annotation of neighboring exons, we filtered out exons sharing their start and end positions with annotated exons, to keep 2,142 exons displaying an unannotated start or end. As many of these had extensive overlap with annotated exons, we further filtered the set to keep 63 exons, 42 of them displaying no overlap with annotated exons, and 21 exons having less than 90% of their sequence overlapping with annotated exons.

### Generating connectivity matrices

We compiled membrane contact and chemical synapse matrices from published electron microscope reconstructions, N2U (Cook et al., 2019; White et al., 1986) and Adults 7 and 8 (Witvliet et al., 2020). Membrane contact data are available for N2U and Adult 8. Chemical synapse data was obtained from three adult animals (N2U, Adult 7 and Adult 8). These sources contain data for each individual neuron (e.g., for each of the six IL2 neurons). Data were summed across the individual neurons corresponding to each neuron type in the single-cell data (e.g., IL2DL, IL2DR, IL2VL, IL2VR were summed for the IL2\_DV class, IL2L and IL2R were summed for IL2\_LR). Only contacts and synapses present across all animals were retained to generate high confidence sets of invariant contacts and synapses.

### Regulatory patterns of neuron transcriptomes

In order to identify distinct regulatory patterns for the transcriptome of each neuron, log-transformed expression values were converted to z-scores from the distribution of expression across all neurons for each gene. A high (low) z-score for a particular gene

in a specific neuron type indicates an upregulated (downregulated) gene relative to the expression in other neurons. For motif discovery in promoters and 3'UTRs, gene z-scores were mapped to their isoform transcripts. Unique isoforms were maintained by applying a simple duplicate removal procedure, which guarantees that no pair of promoters and no pair of 3'UTRs will have a Blast local alignment with E-value  $< 10^{-10}$  (Elemento et al., 2007). For promoter sequences we considered sequences 1KB upstream of the transcriptional start site of each isoform, while for 3'UTRs we considered 1KB from the start of each annotated 3'UTR sequence (or 1KB downstream of the stop codon for transcripts without annotated 3'UTRs). To identify expression patterns of co-regulated transcripts, z-score values across all neuron types were clustered using hierarchical clustering with three different cut-offs (python/scipy fcluster implementation, cosine metric, criterion = 'distance', cophenetic threshold = 1.2, 1.25, 1.37). We chose these thresholds to provide clustering of the data ranging from coarse to fine (16, 48, and 76 transcript clusters). For individual neurons, transcripts were categorized into bins with high to low z-scores based on the distribution of all z-scores across transcripts and neuron types. Z-score bin intervals were defined considering the following percentiles of the overall distribution of z-scores: 2.5%, 5%, 10%, 20%, 80%, 90%, 95%, 97.5%. For each neuron type, the top bin included transcripts with z-scores above the 97.5<sup>th</sup> percentile, the second to top included z-scores between the 95<sup>th</sup> and 97.5<sup>th</sup> percentile, etc. The bottom bin included transcripts with z-scores below the 2.5<sup>th</sup> percentile, the second to bottom included z-scores between the 2.5<sup>th</sup> and 5<sup>th</sup> percentile, etc. To avoid poorly populated bins, any given category containing less than 350 transcripts was merged with the next closest bin toward the center of the distribution.

### Cis-regulatory element discovery

To systematically explore the regulatory effect of short DNA and RNA cis-regulatory elements, we utilized FIRE, a computational framework for *de novo* discovery of linear motifs in DNA and RNA whose presence or absence in a transcript's promoter and 3'UTR regions is informative of regulatory patterns. We ran FIRE in discrete mode including transcript identifiers (Wormbase transcript IDs) along with either their z-score bin categories (for individual neurons) or transcript cluster IDs (for patterns of co-regulated genes). Over representation (yellow) and under representation (blue) patterns are shown for each discovered motif within each category (bin or cluster) of transcripts as well as mutual information (MI) values and z-scores associated with a randomization-based statistical test. All discovered motifs pass a three-fold jackknifing test more than 6 out of 10 times. Each time one-third of the transcripts was randomly removed and the statistical significance of the MI value of the motif was reassessed. For each of the 10 tests, the remaining two-thirds of the transcripts was shuffled 10,000 times and the motif was deemed significant if its MI was greater than all 10,000 MI scores from the randomized sets (Elemento et al., 2007). For every motif identified through FIRE, we defined the regulon for that motif as the collection of transcripts that harbored instances of the motif in their promoters (DNA motifs) or 3'UTRs (RNA motifs).

### Motif families

Motifs with similar nucleotide compositions and regulons were discovered across individual neurons and gene expression patterns. We sought to identify the extent of redundancy between individual motifs and group them into motif families based on their similarity. We included additional motifs in this analysis for known transcription factors (CIS-BP, JASPAR), RNA binding proteins (CISBP-RNA) and miRNA 6-mer seeds (5' extremity of known miRNA sequences of *C. elegans*). To quantify the similarity between nucleotide compositions between motifs we applied TOMTOM (MEME version 5.0.5). For each motif, we used its IUPAC motif sequence to convert it into a MEME formatted motif (iupac2meme function) as input to TOMTOM and compared it against all other discovered and known motifs. We specified a minimum overlap of 5, and an E-value threshold of 10 to identify significant matches. To quantify the extent of overlap between two motif modules, we defined a similarity measure between a module A and B as  $S(A, B) = (G_A \cap G_B) / \min(G_A, G_B)$ , where  $G_K$  is the set of transcripts in module K. We calculated TOMTOM and module similarity scores for all motif pairs. Module similarity scores were deemed significant if  $p < 10^{-4}$  (hypergeometric test). To ensure that motifs are considered redundant only when they are similar both in nucleotide and module composition, we set the module similarity scores to 0 if either the TOMTOM or the module similarity scores were not significant. We clustered the motifs into motif families based on the masked similarity measures of all motif pairs using hierarchical clustering (python/scipy fcluster implementation, cosine metric, criterion = 'distance', cophenetic threshold = 0.9). We set out to identify potential known regulators that represent a given motif family. To this end, we applied TOMTOM to match the motif family members with the binding preferences of known regulators. For each motif family, we counted all the significant TOMTOM scores for every family member compared to a known regulator. We considered a known regulator as a potential match for the motif family, if it had a significant TOMTOM score for more than 2/3 of the family members.

### Associations of motif families and neurons

We set out to assess the regulatory potential of each motif family on each neuron type. Motifs with positive regulatory potential should have consistent patterns across the z-score bins, i.e., predominantly over-represented in genes with high z-scores or under-represented in genes with low z-scores. On the other hand, motifs with negative regulatory potential should be over-represented in genes with low z-scores or under-represented in genes with high z-scores. For each neuron type and each motif, we considered the frequency of transcripts carrying the motif in the top two z-score bins combined ( $f_t$ ), as well as the bottom two z-score bins ( $f_b$ ). To consider a positive association of the motif with the neuron type we required that the motif is: *over-represented in the top two bins* ( $p < 0.005$ ) and *not over-represented in the bottom two bins* ( $p > 0.05$ ), or, *under-represented in the bottom* ( $p < 0.005$ ) *two bins and not under-represented in the top two bins* ( $p > 0.05$ ). To consider a negative association of the motif with the neuron

type we required that the motif is: *over-represented in the bottom two bins* ( $p < 0.005$ ) and *not over-represented in the top two bins* ( $p > 0.05$ ), or, *under-represented in the top two bins* ( $p < 0.005$ ) and *not under-represented in the bottom two bins* ( $p > 0.05$ ). We calculated a  $\text{Log}_2$ -fold ratio ( $\log_2[R] = \log_2\left[\frac{f_t}{f_b}\right]$ ) and an associated p value (hypergeometric test) between the two categories. We reported significant associations ( $|\log_2[R]| > 0.5$  and  $p < 10^{-5}$ ). For each motif family, we report the  $\text{Log}_2$ -fold ratio and signed p value ( $-\text{sgn}(\log_2(R)) * \log_{10}(p)$ ) for the motif member with the lowest p value.

### Cell adhesion molecule by stratum analysis

Given a set of gene expression profiles for the neurons classes in the nerve ring and their memberships in different strata, we can execute standard differential gene expression (DGE) analysis (Soneson and Delorenzi, 2013) to determine which genes are enriched in members of particular strata. Standard DGE analysis involves performing univariate t tests between the gene expression levels of members of a particular stratum versus the members of all the remaining strata. The visual representation of this test can be seen in [Methods S1](#). In detail, the DGE model involves fitting a regression model where the response variables are the gene expression levels for every neuron and the design matrix is a vector of 1 s and  $-1$  s corresponding to the neurons in the two groups that are being compared. The gene expression is logarithm transformed to Gaussianize count-based data (Love et al., 2014). The output of this test is a vector of t-statistics and log-fold changes for every single gene in which this tuple of information can be visualized via volcano plots (Figure S8C). We deem that genes that pass the Bonferroni threshold for multiple comparisons ( $q < 0.05$ ) are significantly enriched or depleted in particular strata.

### Network differential gene expression analysis

Whereas standard DGE analysis is useful for delineating univariate differences between groups of neurons, here we introduce a generalization of DGE, termed “network” DGE (nDGE), to establish the genetic determinants of synaptic formation and maintenance. Unlike DGE where gene expression levels of disjoint groups of neurons are compared, in nDGE, the multiplicative co-expression of genes, between sets of pairs of neurons (representing edges in a network) is compared. The visual representation of the nDGE statistical model can be seen in [Methods S1](#). In nDGE, the response variables are the pairwise co-expression of all genes in all pairs of neurons. On the other hand, the design matrix captures two sets of pairs of neurons, one for each group. Similar to standard DGE, the output of this test is a set of t-statistics and log-fold changes for gene associations. However, unlike standard DGE, the t-statistics and log-fold changes in nDGE capture the effect of co-expression of pairs of genes, one corresponding to the gene observed in the pre-synaptic neuron partner and the other corresponding to the gene observed in the post-synaptic one. To deem a pair of genes significant under nDGE analysis, we also utilize the Bonferroni correction for p values. However, the number of comparisons in nDGE is the square of the number of genes interrogated.

Since nDGE is a generalization of standard DGE, it enables the testing of a variety of hypotheses in addition to what is testable in standard DGE. The types of hypotheses that are tested are encoded in the design matrix of nDGE of which several examples are displayed in [Methods S1](#). [Methods S1](#) shows how standard DGE can be executed through nDGE, by placing 1 s and  $-1$  s in the diagonal of the design matrix corresponding to the neuron groups. Three other types of hypotheses that can be tested are whether particular gene pairs have global effects of synaptic formation across all the neurons, whether there are differential gene co-expression differences in the synapses of two different neurons, or which gene co-expression patterns are implicated in the synapses of an individual neuron. In these scenarios, the design matrix has 1 s where there is a synapse and a  $-1$  where there is membrane contact, but no synapse, restricted to the sets of neurons of interest (all, pair, or one, respectively).

The main caveat in nDGE is the lack of independence of samples that are compared between groups. Since “samples” in nDGE are the co-expression of genes in pairs of neurons, the information from a particular neuron will inevitably be represented multiple times and possibly in different groups e.g., the gene expression from neuron AIA is represented in multiple synaptic gene co-expression values for all synaptic partners of AIA as well as the non-synaptic adjacent partners of AIA (Figure 7B). This lack of independence in the test samples can falsely inflate/deflate the sample variance, which can introduce excess false positives and false negatives. To accurately estimate the null distribution of the nDGE test statistics, we generate randomized “pseudoconnectomes” that respect the topology of the original connectome. Specifically, the pseudoconnectomes preserve the same number of synaptic partners for each neuron and the shuffled synaptic partners are confined to be neurons that have membrane contact (Milo et al., 2003). The latter constraint prevents infeasible pseudoconnectomes where synapses exist between neurons that do not share a membrane contact. Examples of pseudoconnectomes that are generated using the chemical connectome and membrane contact adjacency matrices are displayed in [Methods S1](#). We execute nDGE analysis with the design matrices corresponding to 1000 pseudoconnectomes and compute a t-statistic using the mean and variance of the resulting null distribution.

While the nDGE technique introduced here is a generalization of standard DGE, interrogating the contribution of pairs of genes in the formation and maintenance of synapses between pairs of neurons, nDGE can only account for a single co-expressed gene in either of the two synaptic terminals (pre/post). For this reason, the nDGE model will tend to underestimate the effects of trimer (or higher-order) proteins in the formation and maintenance of synapses. Therefore, it is imperative to keep in mind that lack of significant hits for a particular neuron might not mean that there are no genes implicated in the formation of synapses for that neuron, but rather

that higher-order gene interactions might be at play. Conceptually, it is straightforward to extend the model to higher-order gene interactions, but the prohibitive number of combinatorial gene co-expression enumeration is a computational bottleneck.

Another feature of nDGE is that it is a mass-univariate method, which does not take into account the possibility of interaction of different co-expressed genes in forming or inhibiting synapses. Therefore, the significance results output by nDGE tends to be very conservative with strict control of type 1 errors. This is in contrast with multivariate methods for explaining the genetic bases of connectivity (Kovacs et al., 2020). Due to the relatively high dimensionality of the gene expression data compared to the number of synapses in the chemical connectome, multivariate models tend to overfit and introduce type 1 errors.

## QUANTIFICATION AND STATISTICAL ANALYSIS

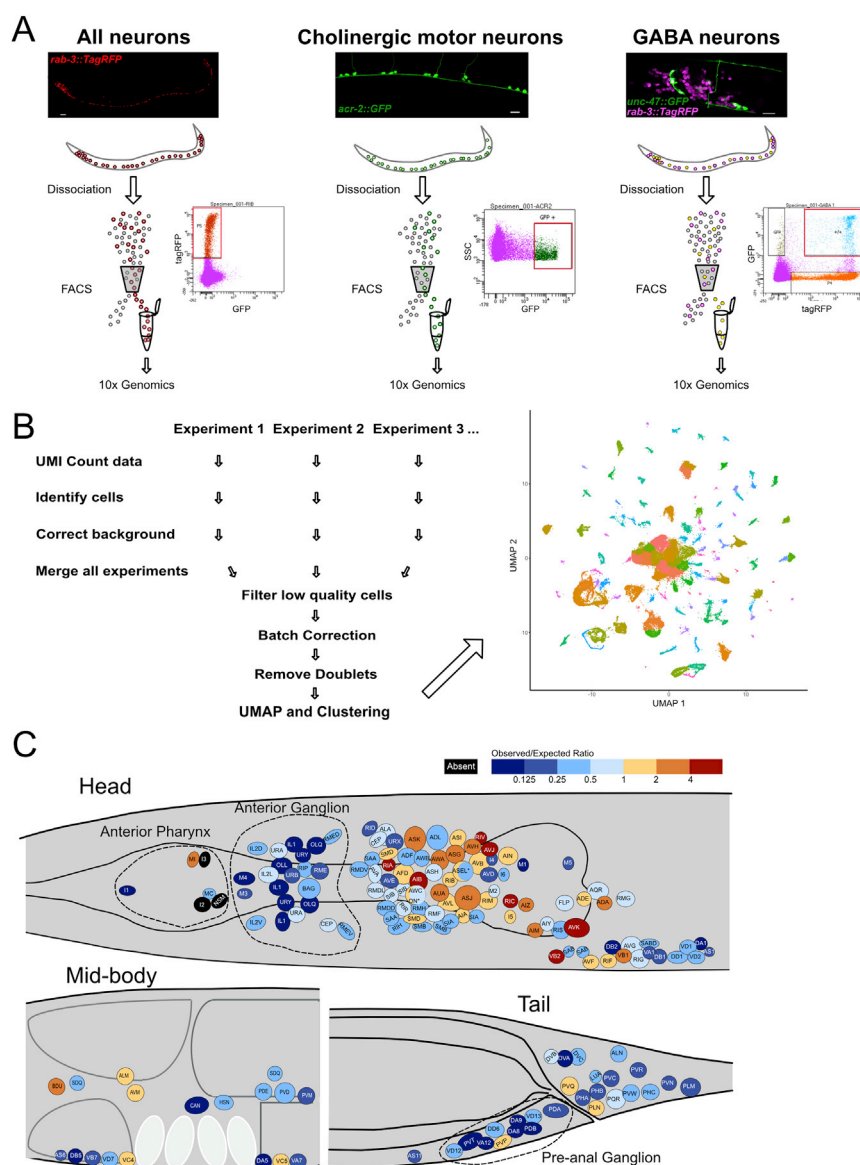
Details of quantification and statistical testing, sample size, center and dispersion are found in the figure legends and [STAR Methods Method details](#) section for individual analyses.

## ADDITIONAL RESOURCES

Data files and information about the CeNGEN Consortium can be found at <https://www.cengen.org>. Single-cell RNA-seq data can be explored, analyzed and downloaded at the CengenAPP, found at <https://cengen.shinyapps.io/CengenApp>.

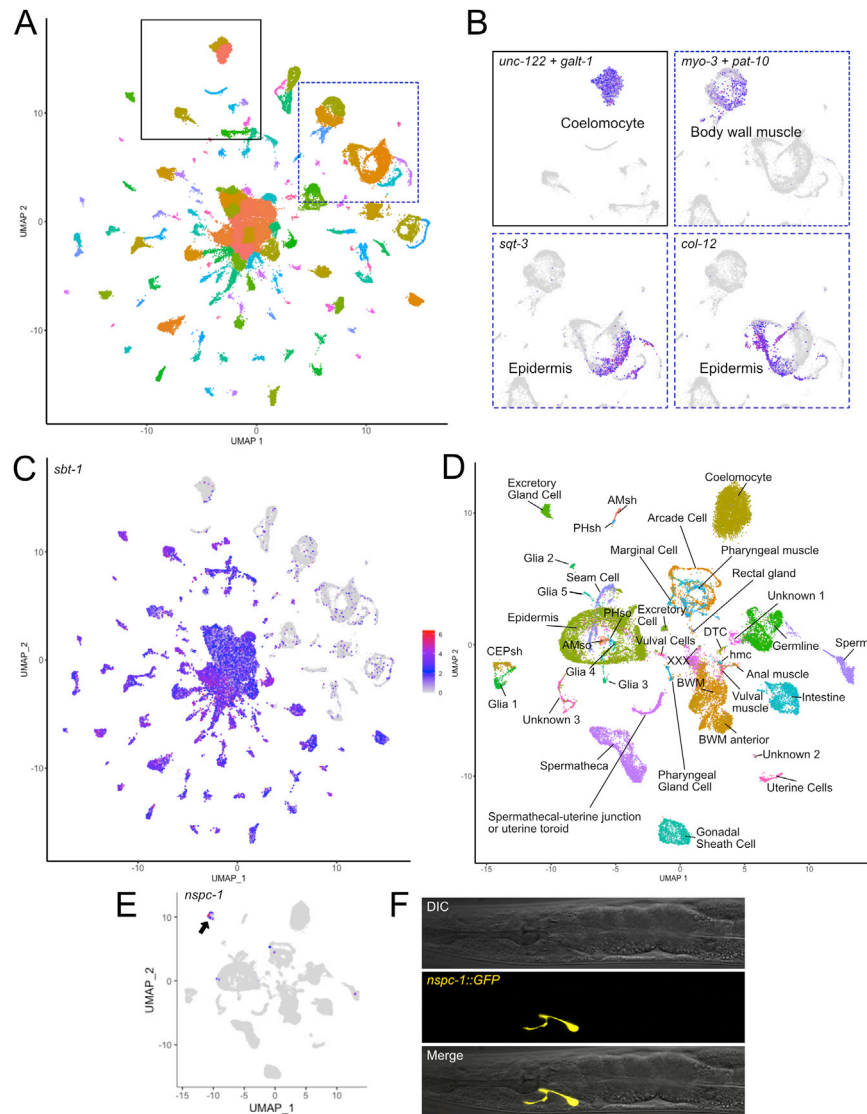


# Supplemental figures



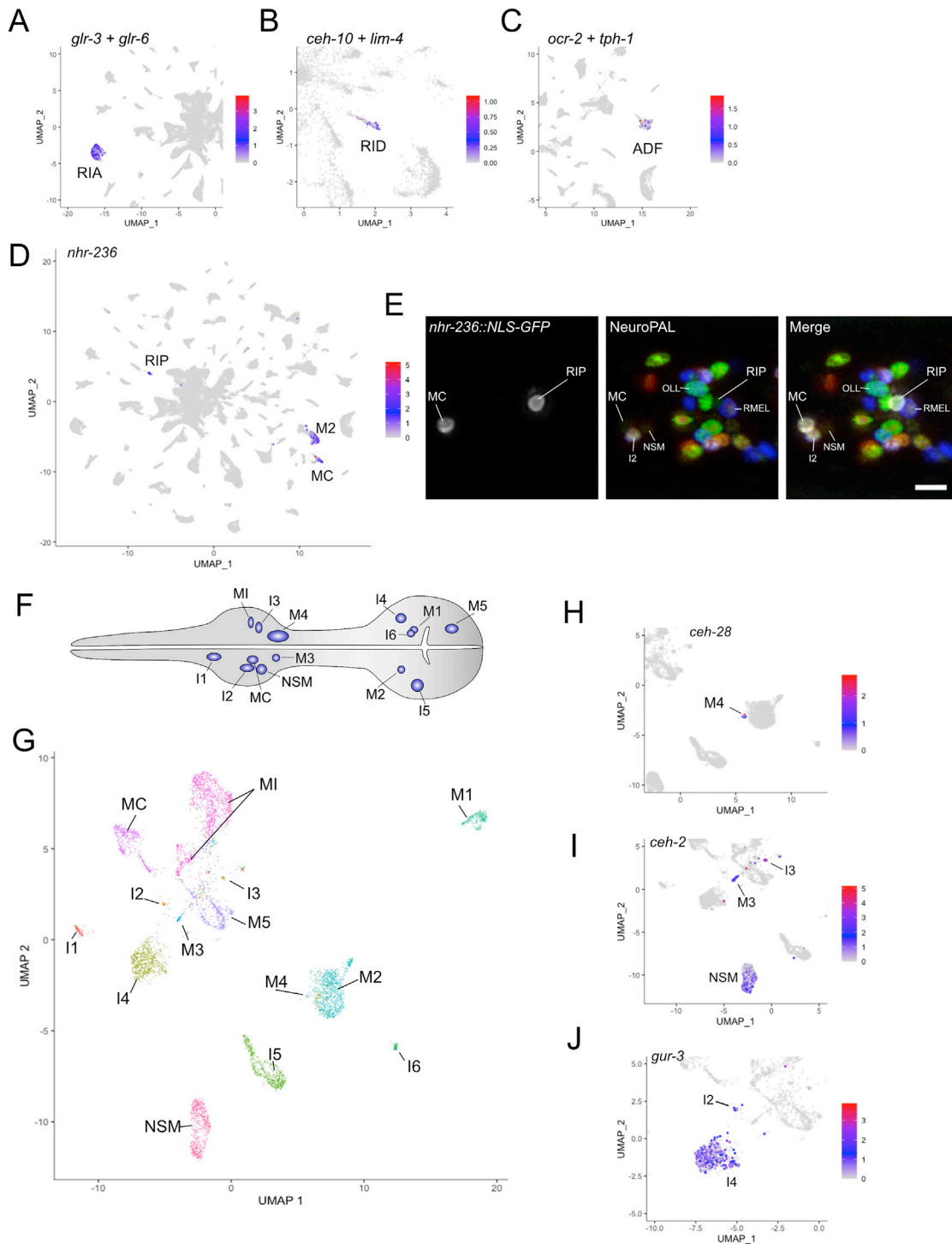
**Figure S1. Isolating L4 larval stage neurons from fluorescent marker strains for scRNA-seq, related to Figures 1 and 2 and Table S1**

(A) Confocal images of the pan-neural marker, *rab-3::TagRFP* (*otIs355*), ventral cord cholinergic motor neurons marked with *acr-2::GFP* (*juls14*) and neurons in the head region dual-labeled with the GABA neuron-specific reporter, *unc-47::GFP* (*oxIs12*) and the pan neural marker strain, *rab-3::TagRFP*. Scale bars = 10  $\mu$ m. L4 animals were treated with SDS-DTT and dissociated with pronase to produce single-cell suspensions. Targeted subgroups of neurons were isolated by Fluorescence Activated Cell Sorting (FACS, red boxes) and collected for scRNA-Seq using the 10x Genomics 3' platform. (B) Results from 17 separate profiling experiments were submitted to a series of processing steps to produce a final merged dataset (see STAR Methods). Right panel shows UMAP projection of all 100,955 cells, colored by cluster. (C) Graphical depiction of relative abundance of each neuron class in cells isolated from the pan-neural marker strain (*rab-3::TagRFP*). The fraction of observed cells for each neuron type was divided by the expected ratio (i.e., # neuron type/302) and annotated for each cell type according to heatmap index. Note under-representation (dark blue/black) of neurons in the anterior pharynx, anterior ganglion and pre-anal ganglion (dashed lines). \* For pairs ASE and AWC, the individual neurons are treated separately.



**Figure S2. Distinguishing non-neuronal versus neuronal cells, related to Figure 1 and Table S1**

A) UMAP projection of 100,955 single cell profiles, colored by cluster. B) UMAP insets from A showing known markers that identify clusters of non-neuronal cells (*unc-122 + galt-1*, coelomocytes; *myo-3 + pat-10*, body wall muscles; *sqt-3* and *col-12*, epidermis). Solid and dashed outlines correspond to similarly boxed sub-regions in A. C) UMAP projection of all cells as in panel A showing the neuropeptide processing gene *sbt-1* is expressed in all neuronal clusters (blue-magenta) and largely absent from non-neuronal clusters (gray). D) Sub-UMAP of all non-neuronal cells, labeled by cell type. E) *nspc-1*, a member of the nematode-specific peptide c (*nspc*) gene family, is restricted to a single non-neuronal cluster (arrow). F) The transcriptional reporter *nspc-1::GFP* is exclusively expressed in the excretory gland cell (yellow). DIC (Differential Interference Contrast). Scale bar = 10  $\mu$ m.



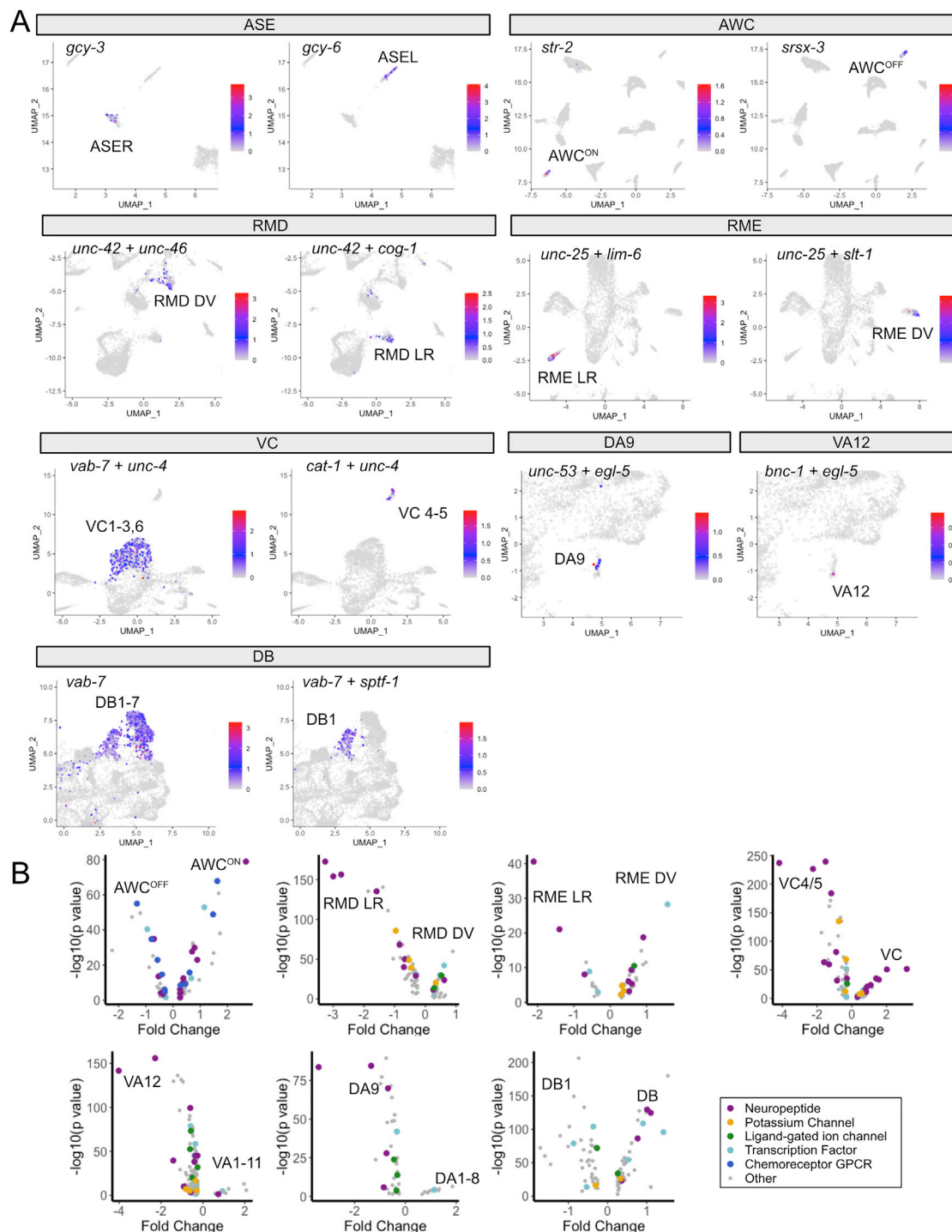
**Figure S3. Annotation of neuronal clusters, related to Figure 1 and Table S1**

A-C) Higher magnification regions of neuronal UMAP from Figure 1B showing selective co-expression of genes used to assign cell identity to individual clusters. A) Co-expression of the glutamate receptor genes *glr-3* and *glr-6* demarcate the RIA cluster. B) Co-expression of the homeodomain transcription factors *ceh-10* and *lim-4* label the RID cluster. C) Co-expression of the transient receptor potential channel (*trp*) gene *ocr-2* and tryptophan hydroxylase *tph-1* label the sensory neuron ADF. D) Neuronal UMAP as in Figure 1B showing normalized expression of the nuclear hormone receptor *nhr-236*. *nhr-236* is primarily detected in three clusters, corresponding to RIP, M2, and MC. E) Z-projection of confocal stack showing *nhr-236::NLS-GFP* expression in RIP and MC in a NeuroPAL strain. M2

(legend continued on next page)

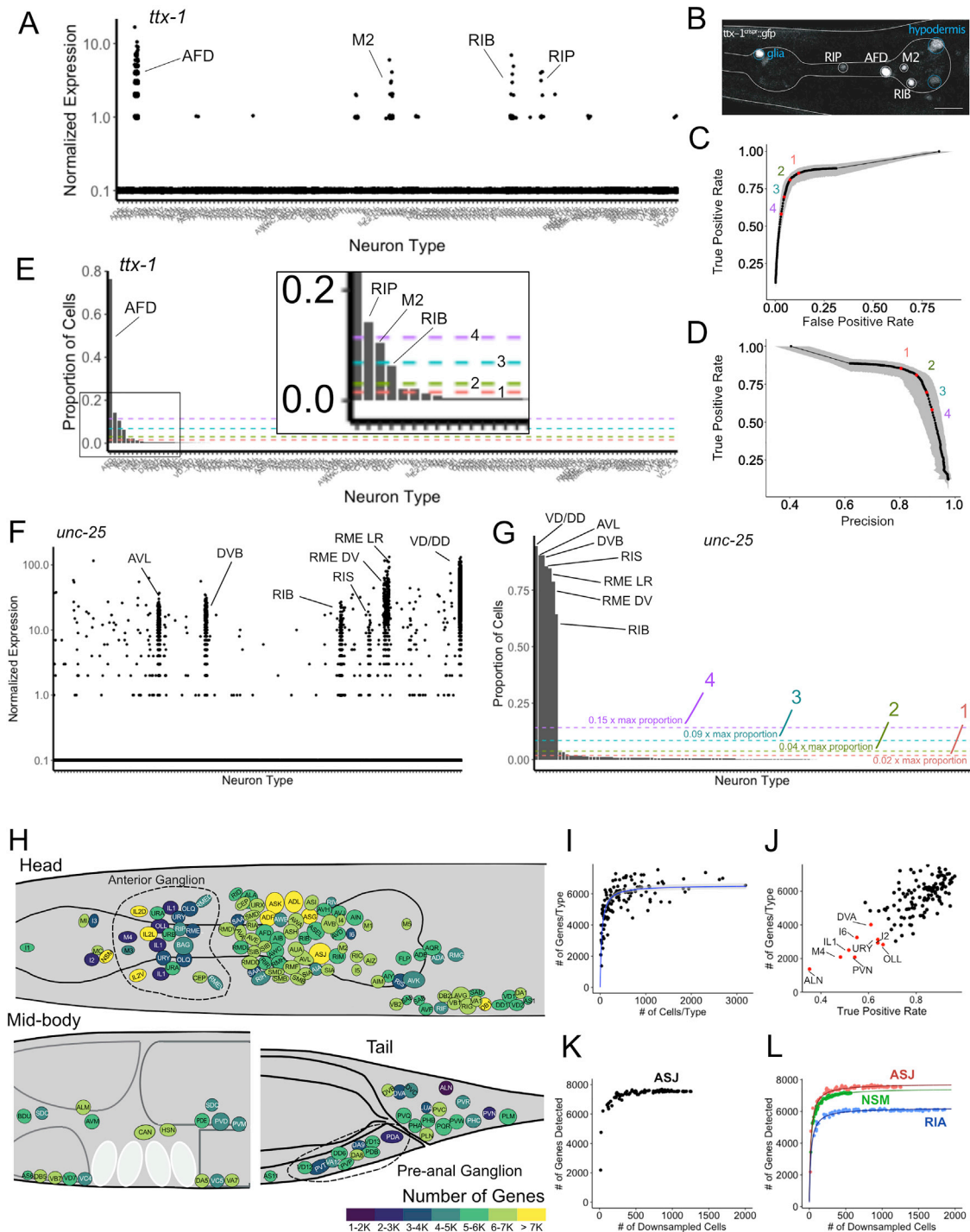
(not shown) also consistently expressed GFP. Scale bar = 5  $\mu$ m. F) Diagram denoting neurons (blue) in pharynx (gray). G) Sub-UMAP of all neuronal clusters expressing the pharyngeal neuron marker *ceh-34* revealed independent clusters for each known pharyngeal neuron type. H-J) Higher magnification regions of pharyngeal UMAP in G showing expression of marker genes used to identify pharyngeal neuron classes. H) The M4-specific homeodomain transcription factor, *ceh-28* is exclusively detected in a small but distinct group of 12 cells. I) *ceh-2*, a marker for I3, M3 and NSM pharyngeal neurons, is restricted to 3 clusters. J) Restricted expression of *gur-3*, a marker for I2 and I4 pharyngeal neurons.





**Figure S4. Identification of neuron subtypes, related to Figures 1 and 2 and Table S1**

A) High magnification regions of neuronal UMAPs (Figure 1B for ASE, AWC, RMD; Figure 1F for VC and RME; Figure 1E for DA, VA, DB) showing expression of known markers for neuron subtypes for ASE (*gcy-3*, ASER and *gcy-6*, ASEL), AWC (*str-2*, AWC<sup>ON</sup> and *srsx-3*, AWC<sup>OFF</sup>), RMD (*unc-42 + unc-46*, RMD DV and *unc-42 + cog-1*, RMD LR), RME (*unc-25 + lim-6*, RME LR and *unc-25 + slt-1*, RME DV), VC (*vab-7 + unc-4*, VC1-3,6 and *cat-1 + unc-4*, VC4-5), DA9 (*unc-53 + egl-5*), VA12 (*bnc-1 + egl-5*) and DB (*vab-7*, DB1-7 and *vab-7 + sptf-1*, DB1). B) Volcano plots of genes that are differentially expressed between neuron subtypes. Inset depicts color coding for selected gene families.



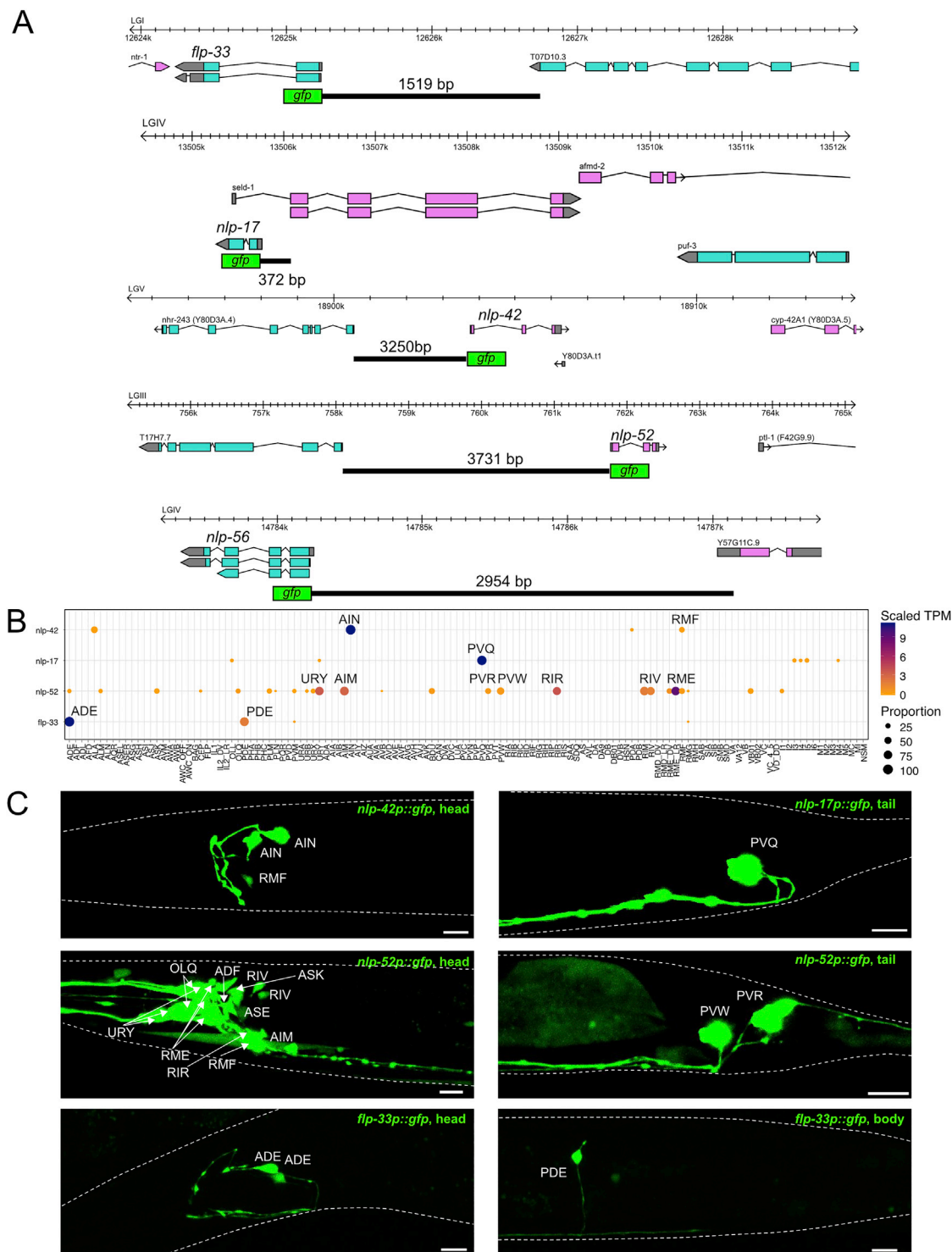
**Figure S5. Establishing expression thresholds, related to Figures 3 and 4 and Table S2**

A) Jitter plot of normalized *ttx-1* expression (y axis) in all neuronal clusters (x axis) shows strongest expression in AFD, M2, RIB and RIP. B) Confocal image of *ttx-1<sup>crispr</sup>::GFP* shows expression in AFD, M2, RIB and RIP neurons and in glia and epidermal cells. C) Receiver-Operator Characteristic (ROC) curve of True Positive Rate (TPR) versus False Positive Rate (FPR) for a range of thresholds (1-4) (red dots) compared to ground truth expression data (see STAR Methods). Increased stringency diminishes both the TPR and FPR. Grey shading represents 95% confidence intervals. D) Thresholds 1-4 (red dots) plotted on Precision-Recall (PR) Curve of Recall (TPR) versus Precision [1 - False Discovery Rate (FDR)]. Grey shading represents 95% confidence intervals. E) The proportion of cells in each neuron-specific cluster expressing *ttx-1*. Inset shows expanded view of boxed region. Thresholds (1-4) are set to different proportions of *ttx-1*-expressing cells in each cluster (see STAR Methods). Note that neurons RIB and M2, that show expression of native *ttx-1<sup>crispr</sup>::GFP*, are excluded by the thresholds 3 and 4. Numbers (1-4) correspond to thresholds in C and D. F) Jitter plot shows strong expression of *unc-25/GAD* in seven known GABAergic neuron types (AVL, DVB, RIB, RIS, RME LR, RME DV, and VD/DD).

(legend continued on next page)

---

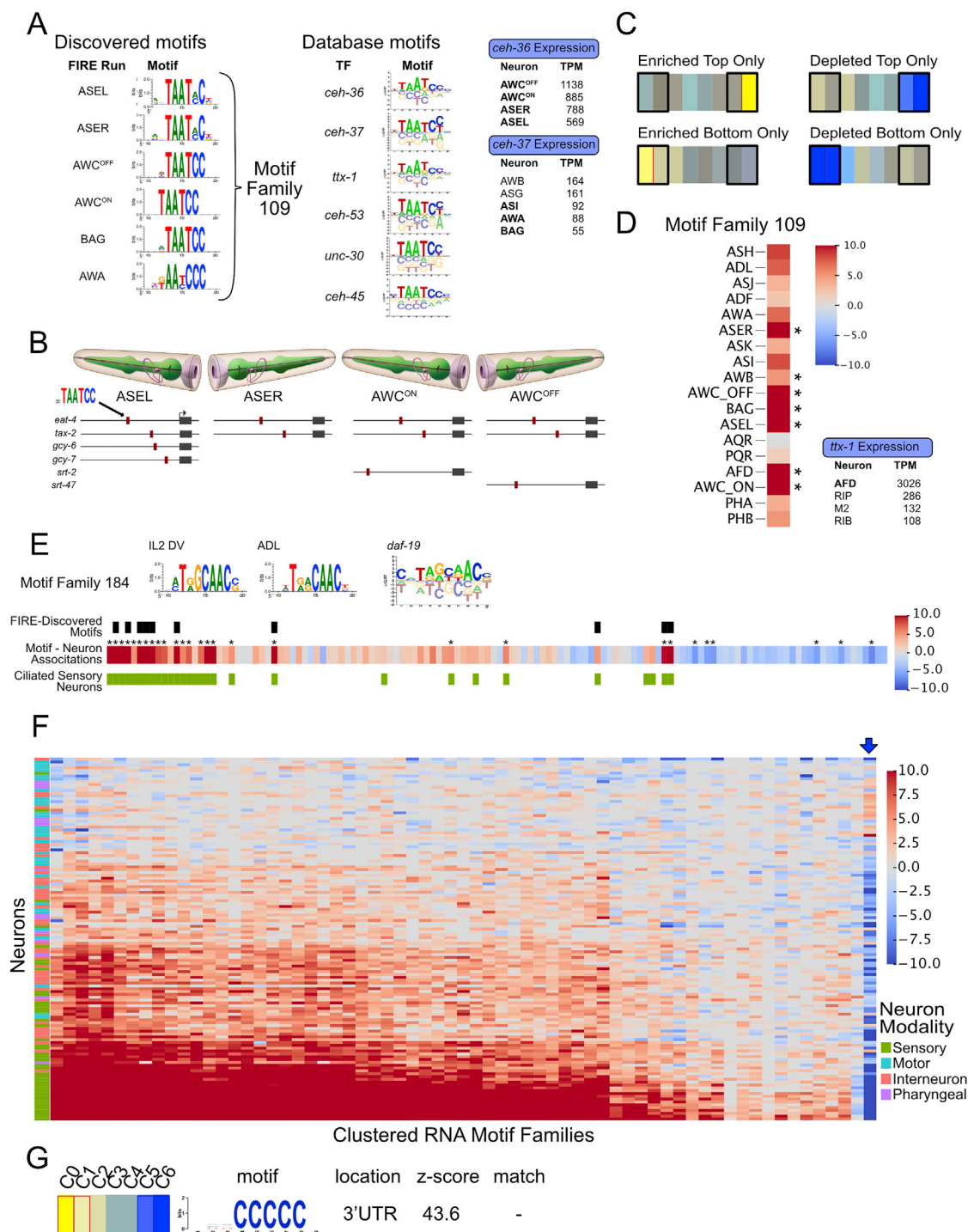
RME, DD, VD) but scattered *unc-25/GAD* expression is also detected in other cell types. G) Bar graph plotting the proportion of cells in each neuron type that express *unc-25/GAD*. Thresholds of increasing stringency (1-4) are set to different proportions of cells in a given cluster that express *unc-25/GAD*. Note that threshold 2 distinguishes known *unc-25/GAD*-positive neurons from other neuron types with lower detected levels of *unc-25/GAD* (see [STAR Methods](#)). H) Cell soma for each neuron type in the head, mid-body and tail regions are colored according to the number of genes detected using threshold 2. Neuron types in the anterior ganglion (dashed line) are among those with the fewest cells and also lower numbers of detected genes. I) Number of genes detected with threshold 2 for each neuron type plotted against the number of cells in each neuron-type cluster. Spearman's rank correlation = 0.783,  $p < 2.2e-16$ . J) Number of genes detected with threshold 2 plotted against the True Positive Rate (TPR) for each neuron type. Spearman's rank correlation = 0.678,  $p < 2.2e-16$ . Neurons with fewest cells, low TPR and number of genes/type are denoted (red). K) Number of genes detected plotted against the number of cells from down-sampling of ASJ cluster. L) Number of genes detected plotted against the number of cells from down-sampling for ASJ (red), NSM (green) and RIA (blue), with corresponding model fits.



**Figure S6. Neuropeptide gene reporter strains validate scRNA-seq data, related to Figure 3**

A) Diagrams depicting upstream regions (black lines) in transcriptional GFP reporter genes. B) Heatmap showing single cell RNA-Seq expression of four selected neuropeptides (*nlp-42*, *nlp-17*, *nlp-52*, *flp-33*) and (C) corresponding confocal micrographs. Neuron types expressing each neuropeptide reporter were determined by co-localization with NeuroPAL markers (not shown). *nlp-42p::gfp* is robustly expressed in AIN and RMF with faint expression in NSM and PHB. *nlp-17* expression at threshold 2 was 480x stronger in PVQ than for six additional neuron types and the *nlp-17p::gfp* reporter was selectively detected in PVQ. GFP reporters for *nlp-52* and *flp-33* (C) are largely congruent with single-cell RNA-Seq results (B). All scale bars = 10  $\mu$ m.

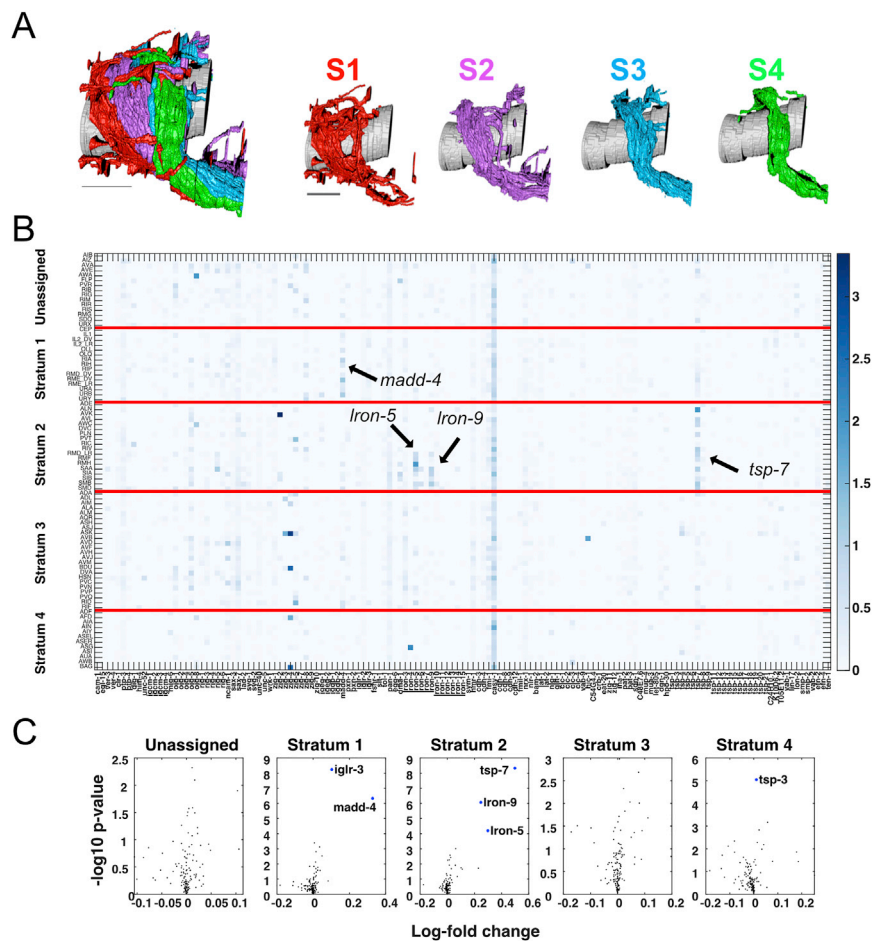




**Figure S7. FIRE discovered known regulatory motifs in individual neurons, related to Figure 6**

A) Left: Motif family 109 includes highly similar motifs discovered in independent FIRE analyses from individual neurons. Center: Motifs from the CIS-BP database for six different K50 homeodomain transcription factors (CEH-36, CEH-37, TTX-1, CEH-53, UNC-30, CEH-45) match motif family 109. Right) Expression (TPM or transcripts per million) of two transcription factors (*ceh-36*, *ceh-37*) in our scRNA-Seq data corresponds to the neurons in which the cognate motif was discovered. B) The TAATCC motif (red boxes) in the 5' regions of genes enriched in ASEL, ASER, AWC<sup>ON</sup> or AWC<sup>OFF</sup>, including genes expressed in all four neurons (top rows) as well as genes expressed in subsets of these neurons. Expression of *eat-4*, *tax-2*, *gcy-6*, *gcy-7* and *srt-2* is *ceh-36*-dependent in these neurons. C) Schematic showing the z-score bins tested for relative enrichment for motif-neuron associations (see STAR Methods). D) Motif-neuron associations (from FIRE Main Figure B for motif family 109 in a subset of neurons) showing correspondence to expression of K50 homeodomain proteins (black arrows). Colors indicate  $-\log_{10}(p \text{ value})$  for positive associations and  $\log_{10}(p \text{ value})$  for negative associations. Asterisks denote significant associations ( $p \text{ val} < 1e-5$ , log fold (legend continued on next page)

change > 0.5). E) Top: Two similar motifs discovered in FIRE runs of individual neurons IL2\_DV and ADL are members of motif family 184 that matches the DAF-19 motif from CIS-BP. Other members of the DAF-19 motif family were discovered in FIRE analysis of 8 additional individual neurons (black boxes). Bottom: Motif family 184 showed significant positive associations (asterisks) with 22 ciliated sensory neuron types (green boxes) and significant negative associations with some motor neurons (blue boxes with asterisks). F) Heatmap of log-transformed p values showing motif-neuron associations for RNA motif families (columns) across all neuron types (rows). Note the strong positive signal in sensory neurons, and the one motif family with largely negative associations with neurons (blue arrow). G) Representative RNA motif (from ADL FIRE run) of family 23 showing a poly-C RNA sequence in the 3' UTR.



**Figure S8. Differential expression of CAMs between strata, related to Figure 7**

A) Cartoon representation of nerve ring strata. Adapted with permission from Nature, Structural and developmental principles of neuropil assembly in *C. elegans* (Moyle et al., 2021). B) Heatmap of CAM expression (columns) in neurons grouped by strata (rows). Black arrows indicate four significantly enriched CAMs with the highest log fold changes between strata. Colors indicate log-transformed expression values. Red lines separate strata. C) Volcano plots showing log-fold change (x axis) by  $-\log_{10} p$  value (y axis) for each CAM in the neurons in each stratum compared to all other neurons. The six labeled genes were significantly enriched in the strata shown.