

Supplementary Materials for

Cell type profiling in salamanders identifies innovations in vertebrate forebrain evolution

Authors: Jamie Woych¹†, Alonso Ortega Gurrola^{1,2}†, Astrid Deryckere¹†, Eliza C. B. Jaeger¹†, Elias Gumnit¹†, Gianluca Merello¹, Jiacheng Gu¹, Alberto Joven Araus³, Nicholas D. Leigh⁴, Maximina Yun^{5,6}, Andrés Simon³, Maria Antonietta Tosches^{1*}

Affiliations:

¹Department of Biological Sciences, Columbia University; New York City, 10027 New York, USA

²Department of Neuroscience, Columbia University; New York City, 10027 New York, USA

³Department of Cell and Molecular Biology, Karolinska Institute; Stockholm, Sweden

⁴Molecular Medicine and Gene Therapy, Wallenberg Centre for Molecular Medicine, Lund Stem Cell Center; Lund University, Sweden

⁵Technische Universität Dresden, CRTD/Center for Regenerative Therapies Dresden; Dresden, Germany

⁶Max Planck Institute for Molecular Cell Biology and Genetics; Dresden, Germany

*Corresponding author. Email: mt3353@columbia.edu

†These authors contributed equally to this work.

This PDF file includes:

Materials and Methods

Supplementary Text

Figs. S1 to S19

Captions for Movies S1 to S3

Caption for Data S1

Other Supplementary Materials for this manuscript include the following:

Movies S1 to S3

Data S1

Materials and Methods

Animals

Adult *Pleurodeles waltl* were obtained from breeding colonies established at Columbia University and Karolinska Institute. Animals were maintained in an aquatics facility at 20°C under a 12L:12D cycle (63). All experiments were conducted in accordance with the NIH guidelines and with the approval of the Columbia University Institutional Animal Care and Use Committee (IACUC protocol AC-AABF2564). Experiments were performed with adult (5-19 months) male and female salamanders, and stage 36, stage 41, stage 46 and stage 50 embryos and larvae (staged according to the Gallien and Durocher 1957 atlas (64)).

Brain dissociation and single-cell capture

For brain dissociation, animals were deeply anesthetized by submersion in 0.2% MS-222. Animals were perfused transcardially with ice-cold oxygenated Amphibian Ringer's solution (96 mM NaCl; 20 mM NaHCO₃; 2 mM KCl; 10 mM HEPES; 11 mM glucose; 2 mM CaCl₂; 0.5 mM MgCl₂), and then decapitated. The whole brain or telencephalon was dissected out, and embedded in 4% LM Agarose in Ringer. The brain was then sliced coronally on a vibratome into 500 μm sections, and cut into 500 μm cubes in cold carbogenated Ringer. The tissue pieces were transferred to a 5 mL tube in 2.5 mL Dissociation Buffer (20 U/mL Papain, 200 U/mL DNase, 5 μg/mL Liberase, 1 μM TTX), and incubated at RT for 30 minutes on a rotator. Following enzymatic dissociation, the tissue was mechanically dissociated while on ice by trituration with fire-polished, silanized glass pipettes of decreasing tip diameter. The supernatant was passed through a 100 μm cell strainer, with fresh carbogenated calcium-free Hibernate A media (BrainBits) added between each subsequent pipette until a uniform suspension was obtained. Calcium-free Hibernate A media was added up to a total volume of 20 mL, and the solution was passed through a 70 μm cell strainer. To the bottom of the tube, 5 mL of 4% BSA in calcium-free Hibernate A with phenol red was added. To filter out cell debris, the cell suspension was centrifuged through the density gradient at 300xg (4°C) for 5 min. The supernatant was removed, and the pellet resuspended in 50-100 μL calcium and magnesium-free Hibernate A media (BrainBits). The cell concentration was determined by counting on a Fuchs- Rosenthal chamber with trypan blue, and diluted to 1000 cells/μL in calcium and magnesium-free Hibernate A.

The cells were loaded into a 10x Chromium Chip G with a targeted cell recovery of 5000-8000 cells for GEM Generation and cell barcoding. Single-cell RNA-seq libraries were prepared with 10x Chromium Next GEM Single Cell 3' Reagent Kit v3.0 or v3.1 (Dual Index).

Developmental Single Cell Collection: For the developmental scRNA-seq dataset, telencephali from 35 larvae at stage 36, a total of 26 larvae at stage 41, 10 larvae at stage 46, and 5 larvae at stage 50 were isolated. Tissue pieces were dissociated for 30 min at room temperature using the Papain dissociation kit (Worthington) in calcium-free Hibernate A (BrainBits). After processing and concentration, cells were resuspended in calcium and magnesium-free Hibernate A (BrainBits) at a concentration of 1000 cells/μL. Cells were counted in a haemocytometer chamber and immediately processed for single-cell GEM formation (10x Genomics, single cell RNA sequencing 3', Chromium V3.1). Cells from stage 36, 41 and 46 larvae were resuspended at ~1500 cells/μL and multiplexed using 4 tags from the Chromium Next GEM Single Cell 3' v3.1 prior to GEM formation.

Analysis of single-cell RNA sequencing data

Reference transcriptome

Tissue was harvested from the brain of a female adult *Pleurodeles waltl*. Tissue was immediately frozen on dry ice and stored at -80°C. Approximately 25mg of brain tissue was then used for tissue pulverization in liquid nitrogen. Tissue powder was used for RNA extraction using Total RNA Purification Kit (Cat. 17200 from Norgen Biotek) as per manufacturer's recommendation. Genomic DNA was removed via on column DNA removal using Norgen RNase free DNAase kit (Cat. 25710). Input QC of the RNA was performed on the Agilent Bioanalyzer instrument, using the Eukaryote Total RNA Nano kit to evaluate RIN and concentration. RIN value obtained was 9.6. The sample libraries were prepared according to Pacbio's Procedure & Checklist – Iso-Seq™ Express Template Preparation for Sequel® and Sequel II Systems, PN 101-763-800 Version 02 (October 2019) using the NEBNext® Single Cell/Low Input cDNA Synthesis & Amplification Module, the Iso-Seq Express Oligo Kit, ProNex beads and the SMRTbell Express Template Prep Kit 2.0. 300 ng RNA was used as input material. The samples were amplified 12 cycles. In the purification of amplified cDNA the standard workflow was applied (sample is composed primarily of transcripts centered around 2 kb). Quality control of the SMRTbell libraries was performed with the Qubit dsDNA HS kit and the Agilent Bioanalyzer High Sensitivity kit. Primer annealing and polymerase binding was performed using the Sequel II binding kit 2.0. Libraries were sequenced on the Sequel IIe instrument, using the Sequel II sequencing plate 2.0 and the Sequel® II SMRT® Cell 8M, movie time 24 hours and pre-extension time 2 hours. (1 SMRT cell per sample.) Raw sequencing data was then fed into IsoSeq3 (<https://github.com/PacificBiosciences/IsoSeq>), as per standard instructions, to generate fasta outputs of high quality and low quality reads. The high quality fasta output was taken for further use as a reference transcriptome.

This IsoSeq transcriptome and previously published short-read based *de novo* transcriptome (iNewt, ref (65)) were used to generate a combined reference transcriptome. TransDecoder was used to generate predicted peptide sequences for each IsoSeq transcript (66). For transcripts with multiple predicted peptide sequences, only the sequence with the lowest e-value was kept. These predicted peptides were provided to EggNOG-mapper (v2.16), along with the published iNewt peptide sequence to generate a transcript to gene name mapping (67). All parameters were set to default (Taxonomic Scope = Auto, Orthology Restrictions = Transfer annotations from any orthologue, Gene Ontology Evidence= Transfer non-electronic annotations, PFAM refinement = Report PFAM Domains from Orthologues, SMART annotation = Skip). Peptide sequences that did not receive a “preferred name” from the EggNOG-mapper were removed. This would occur when either the sequence was not assigned an orthologous sequence, or the assigned sequence was not characterized enough to receive a common gene name. The name for the RNA transcript sequence that correlated with the peptide sequence used for EggNOG mapping was matched with the corresponding “preferred name” from the EggNOG output to generate a tgMAP for subsequent Alevin analysis described below (68). The merged transcriptome was then filtered to only contain genes which were included in this tgMAP file (again, based on their assignment of a “preferred name” from EggNOG). This final merged reference transcriptome contained 194,403 transcripts of 17,111 uniquely annotated genes.

Read alignment

Reads were assigned to the combined *Pleurodeles waltl* transcriptome using Alevin (Salmon v1.6.0) (68). Library type was set to automatic and keepCBFraction was set to 1. All

other parameters were set as default. The stage 36, 41 and 46 libraries were multiplexed (as described above) and required demultiplexing prior to analysis. For this, multiplex tag reads were also assigned using Alevin.

Data QC and filtering

For adults: Count matrices were obtained after alignment with Alevin (Salmon v1.6.0) and knee plots were generated for each library. Libraries were processed independently to filter out droplets with low UMI counts, according to the plot inflection point. These filtered count matrices were used as input for the R package Seurat 4.0.2 (69), to generate individual Seurat objects. All objects were merged together into a single Seurat object containing all libraries, which was then filtered to keep cells with the following quality statistics: percentage mitochondrial genes <15%, number of genes/cell >800; this produced a dataset of 52,638 cells. After Louvain clustering, low-quality clusters were identified with the following criteria: number of genes per cell below the average of the dataset, percentage of mitochondrial genes above dataset average, and absence of cluster-specific marker genes (these low-quality clusters expressed high levels of housekeeping, mitochondrial, and ribosomal genes). Ten percent of cells from these low-quality clusters, and 10% from the rest of the dataset were used to train a Support Vector Machine (SVM) classifier (7). The SVM classifier identified other low-quality cells in the entire dataset, which were filtered out to obtain a final Seurat object of 36,116 cells. Additional quality control of this final dataset was performed by analyzing the number of genes per cell, number of UMIs per cell, and percentage of mitochondrial gene expression both by library source and by animal source.

For Development: The same protocol was used as for the adult brain, except that cells from stage 41 and 50 telencephalon samples were initially filtered at 1000 reads and 15% reads from mitochondrial genes, before being cleaned by SVM as described above. Additional cells from the stage 36, 41 and 46 telencephalon were initially filtered at 700 reads before demultiplexing using the Seurat package. Only cells labeled as singlets were kept for further analysis, and were cleaned by SVM before merging with the un-multiplexed stage 41 and 50 samples.

Clustering

The adult 36,116-cell dataset was clustered and analyzed using the R package Seurat 4.0.2 (69). Raw data was normalized using the SCTransform function from Seurat, regressing the data by the following variables: animal, percentage of mitochondrial gene expression, number of genes per cell and number of UMIs per cell. The top 2000 variable genes were used for downstream analysis. After principal component analysis (PCA), the first 50 principal components were used for clustering and UMAP embedding. For clustering, we used the default clustering method and the original Louvain algorithm with $res = 2$; clusters were annotated based on the expression of marker genes. Neurons were identified as clusters expressing neuronal markers such as *Syt1*, *Snap25*, *Slc17a7*, *Slc17a6*, *Gad1* and *Gad2*. To identify neuronal clusters with higher resolution, we filtered the original count matrix to keep only neurons, and used this matrix for clustering and UMAP embedding, as described above (variables regressed: animal, percentage of mitochondrial gene expression, number of UMIs per cell and number of genes per cell). The top 2000 variable genes were used for PCA, and the first 180 PCs for Louvain clustering with resolution = 6. At this resolution, some parts of the dataset were overclustered and some others underclustered, as revealed by the expression of marker genes. These clusters

were further analyzed and merged or split on the basis of differential expression of marker genes. Validation of clustering resolution was made by inspecting the presence of cluster-specific marker genes in the salamander neuronal dataset (Fig. S3). The exact steps carried out for manual curation of neuronal clusters can be found in our open repository (<https://github.com/ToschesMA/Salamander-telencephalon>). After these curations steps, the final clusters were renamed according to their excitatory (*Slc17a7*, *Slc17a6*) or inhibitory (*Gad1*, *Gad2*) identity, as well as their anatomical location (telencephalon or diencephalon): TEGLU refers to telencephalic glutamatergic neurons, TEGABA to telencephalic GABAergic neurons, and DIME to diencephalic and mesencephalic neurons. This final neuronal dataset included 29,294 cells and 114 clusters.

Developmental data: SVM cleaned data from stage 36, 41, 46 and 50 was combined into one Seurat object using the “merge” function, which contained 23,100 cells over 39 clusters. Cell cycle scores were computed using the CellCycleScoring function based on Seurat’s built in list of cell cycle genes. Non-neuronal clusters were identified and then removed on the basis of marker genes. The remaining 22,836 neurons were normalized using the SCTransform function while regressing for RNA count, stage, percent mitochondrial genes, G1 score, and G2M score. The top 3000 variable genes were used for downstream analysis, and the first 35 PCs were used for clustering analysis and UMAP embedding; Louvain clustering with a clustering resolution of 3 yielded 53 clusters. This sample was further filtered for trajectory inference by removing non-telencephalic clusters, identified for lack of expression of the telencephalic marker *Foxg1* (Fig. S9B) (70). The remaining 20,261 telencephalic cells were normalized as above and clustered (35 PCs, resolution 1.5) 1.5 to generate 30 clusters. Label transfer against the adult data: Label transfer within the Seurat package was performed to annotate the developmental clusters based on differentiated adult cell types. This algorithm identifies transfer anchors between two data sets (in this case the cleaned adult data set, with only neurons and ependymolgia that received a defined annotation, and the above described telencephalic developmental data) to allow a comparison of cell type identity between the two, meaning that developmental cells can be classified based on which adult cell they are most similar to. The function FindTransferAnchors was run using `dims= 1:50` and `k.anchor=10`, and the function TransferData was run using `k.weight=20`.

Hierarchical clustering of adult neuron clusters

A matrix of average gene expression profiles by cluster was calculated for the neuronal object using the AverageGeneExpression function in Seurat on the SCTransform normalized neuronal dataset. After subsetting this matrix to keep the top 2000 variable genes, distances were calculated as $1 - \text{cor}(x)$, where cor is Spearman correlation. From this distance matrix, hierarchical clustering was performed using Ward’s minimum variance method, squaring dissimilarities before clustering (‘ward.D2’).

Trajectory inference

The developmental telencephalic clusters were annotated according to their corresponding pallial region as described above. *Neurog2* expression was used to identify pallial progenitors, as this gene is expressed predominantly in committed neural progenitors in vertebrates (71, 72). We then included clusters that corresponded to DP, MP, LP,VP, MT, and AMY, as well as their progenitors for further analysis. This data was SCT normalized and regressed as described above for the other developmental data and clustered using 53 dimensions with a clustering resolution

of 1. This new UMAP underwent label transfer against only the adult neurons using the same parameters as above. Cells that scored above 0.75 predicted.id for a pallial cluster (i.e. DP, VP, AMY, etc.) were re-assigned to that cluster identity. Cells in terminal clusters that did not meet this threshold were removed as completely overlapping clusters confounds trajectory analysis. The subset was provided to Slingshot (v2.5.1) (32) for trajectory analysis using default parameters with extend = "n", and pseudotime was calculated along each of the inferred trajectories by measuring the distance of each cell along a Slingshot-defined trajectory from a common progenitor population to one of seven terminal branches. Cells received a "curve weight" representing their assignment to a given principle curve trajectory. Genes differentially expressed between the dorsal and ventral trajectories were determined using the FindMarkers function within Seurat (default parameters - gene.use set to all, and test.use set to bimod) between cells that received above a 0.95 curve weight and pseudotime value above 7 for each trajectory. Genes were selected for heatmap plots in Fig. 3 if they showed greater than 10% difference in the absolute proportion of cells expressing these genes between clusters, and 2 fold greater percentage of cells that express the gene between the clusters that compose each trajectory. p values for all genes identified were below 1e-9. Genes were then ordered on the heatmap based on when they reached peak expression along pseudotime. A heatmap was generated using the ComplexHeatmap R package (73), and heatmap colors were assigned based on the root mean square of cell expression levels for each gene using the scale function.

Comparison with mouse in situ hybridization from the Allen Brain Atlas

The comparison of *Pleurodeles* scRNAseq data with *in situ* hybridization data from the Allen Adult Mouse Brain Atlas (ABA) was performed as described in Colquitt et al. 2021 (8). Briefly, digitized ABA *in situ* hybridization data were downloaded using the R package *cocoframer* (<https://github.com/AllenInstitute/cocoframer>). This dataset consists of average pixel intensity values for every gene and every brain region annotated in the ABA. The dataset was filtered to include only brain regions annotated as "cortex" (CTX) and "cerebral nuclei" (CNU); CTX includes all mouse pallial regions, including the pallial amygdala; CNU includes mouse subpallial regions, such as the striatum and the septum and their subdivisions. The top 1000 high variable genes in the *Pleurodeles* telencephalic neurons dataset were selected, and intersected with the list of genes available in the ABA dataset, yielding 496 genes. These genes were used to calculate expression correlations between mouse brain regions and *Pleurodeles* clusters, as described in Colquitt et al. 2021 (8).

Cross-species comparisons of transcriptomics data

One-to-one orthologues were identified using EggNOG mapper (67), run using the same parameters described above except now with taxonomic scope = vertebrates, and orthology restrictions = One to one orthologues. EggNOG orthology assignments for lizard (*Pogona vitticeps*) and turtle (*Chrysemys picta*) were taken from Tosches et al. 2018 (7). The predicted names column in this output was then compared to similar EggNOG mapper results performed on other species in order to determine which genes were to be included for cross species comparison. Only one-to-one orthologs in all species analyzed were used for cross-species comparisons (with the exception of SAMap, see below).

We used the R package Seurat, as well as other manifold integration algorithms, to compare scRNAseq data across species.

Seurat: Using Seurat's integration pipeline (31), we generated two integrated datasets. The first integrated dataset (Fig. 4) included scRNAseq data from the lizard *Pogona vitticeps* (35), the turtle *Trachemys scripta* (7) and from salamander (this study). For consistent comparisons, we subsetted the original datasets in order to only include cells that were sampled from equivalent brain regions. The lizard dataset from Hain et al. (35) includes and extends the initial lizard dataset by Norimoto et al. (34) and comprises cells from the adult telencephalon, diencephalon and midbrain. Using the brain region annotation provided by Hain et al., we selected neurons annotated as telencephalic, and among those, only neurons collected and sequenced with the 10x Genomics v3 kit. Next, we annotated the Hain et al. telencephalic dataset on the basis of the original identities assigned by Norimoto et al. (34) (Fig. S10). For the turtle dataset, we included excitatory ("e") and inhibitory ("i") clusters, which included interneurons and some subpallium (7), and excluded the unidentified clusters. For salamanders, we included cortical pallium, amygdala, septum, striatum and interneurons, excluding cells from the olfactory bulb, because the olfactory bulb was not sampled in the reptilian datasets. After subsetting, each dataset was normalized independently using Seurat's SCTransform v2 function, which corrects for differences in sequencing depth (74). Each dataset was regressed by percent of mitochondrial genes and animal of origin. These three datasets were brought together in a list class object, from which the integration features were calculated. Using the SelectIntegrationFeature function in Seurat, we selected the top 2000 genes that are variable across datasets and prepared the dataset for integration using these 2000 genes. Pairs of mutual nearest neighbors (anchors) were identified using the FindIntegrationAnchors function with the arguments reduction="cca" (Canonical Component Analysis) and normalization.method = SCT. Integration was carried out with the Seurat function IntegrateData, normalizing with SCT, using the anchor sets obtained from FindIntegrationAnchors and 80 dimensions. Downstream processing was performed as described above, calculating PCA with 200 PCs and UMAP with 80 dimensions. Clustering analysis was performed with the default method, using the SLM algorithm and clustering res = 2.1. The resulting integrated object consisted of 39,391 cells from the three species, with a total of 65 clusters from which 27 are glutamatergic clusters and 38 GABAergic clusters. We carried out additional analysis with the same datasets and pipeline, but changing the number of genes for integration (500, 1000, and 3000) while keeping the same number of dimensions (80) as a validation of the robustness of the integration analysis (Figure S12).

The second integrated dataset included cells from lizard (Hain et al. 2022, (35)), turtle (Tosches et al. 2018 (7)), mouse (Yao et al. 2021 (43)), and salamander (this study). The goal of this integration was to compare neurons from the medial and dorsal pallia of tetrapods. Therefore, we subsetted these datasets in order to include only cells from the dorsal and medial pallium (scRNAseq data from the mouse ventrolateral pallium derivatives are not available). For lizards, we subsetted the dataset to keep only medial cortex (MCtx), dorsal cortex (DCtx), MGE-derived interneurons and CGE-derived interneurons. Similarly, the turtle dataset was subsetted to keep cells from the dorsal cortex (DC), medial cortex (MC), dorsomedial cortex (DMC), and interneurons. (Note on nomenclature: the reptilian cortex includes two hippocampal regions, called medial and dorsomedial cortex. In lizard, clusters from the hippocampus could be identified, but were not mapped precisely to hippocampal subdivisions; therefore we use the term MCtx to annotate cells that belong to brain areas traditionally called medial and dorsomedial cortex). For salamanders, we included medial pallium, dorsal pallium, and interneurons. The Yao et al. mouse dataset was subsampled to keep only 25,000 cells from the original dataset; this dataset includes neurons from all derivatives of the medial and dorsal pallia, including

hippocampus, neocortex/mesocortex, subiculum, entorhinal and cortical interneurons. Following an approach similar to what is described above, all datasets were independently normalized using SCT v2. For turtle, lizard and salamander, the same variables to regress were used. For mouse normalization, the only variable to regress was `external_donor_name_label`. To select genes for integration, we used a different approach given that the greater diversity of mammalian cells biases the integrated features towards mammalian genes. For balancing the list of integration genes, we calculated the top 5000 variable genes from each individual dataset and took the intersection of these lists. Additionally, we removed from the integration list any gene that had a ‘salt-and-pepper’ expression, to keep only genes that were expressed in at least 20% of cells of at least one cluster and one species. After these steps, 877 genes were left for integration. The rest of the analysis was carried out as described above, using a clustering resolution of 1.2. The final integrated dataset for all four species consisted of 39,799 cells and 52 clusters.

To identify high-level molecular similarities among different species’ cells and integrated clusters, we performed a cross-species taxonomy analysis for both integrated datasets, using the R `speciesTree` package developed by Bakken et al. 2021 (<https://github.com/huqiwen0313/speciesTree>). Average expression data for each integrated cluster were calculated in the integration space (`assay="integrated"`) using the Seurat function `AverageGeneExpression`. Distances of integrated clusters were calculated as $1 - \text{cor}(x)$ (Spearman correlation method) and hierarchical clustering of this distance matrix was calculated with Ward’s method. Branches of the resulting dendrogram were color-coded according to the proportion of each species’ cells in the integrated cluster, calculated on the basis of the entropy of the normalized cell distribution, as described in Bakken et al. 2021. The same species mixing color scheme was used in the UMAP plots of Figures 4A and 5B, where each cluster is represented by a circle of size proportional to the number of cells in the cluster, and colored according to species mixing.

Other integration algorithms: the same scRNAseq datasets described above were integrated with alternative approaches. For integration with the R package `Harmony` (75), the function `RunHarmony` was called from within Seurat, with the parameter `assay.use = "SCT"`. The first 80 dimensions were used to compute the nearest neighbors graph and the UMAP embedding.

Using the python package `scVI-tools` (single-cell Variational Inference (76)), we performed the integration using the same genes that were used for Seurat integration and the raw counts for the individual objects described above in a single merged Seurat object using Python from R with the R `reticulate` package. (2000 genes for the analysis in Fig. 4 and 877 genes for the analysis in fig. 5). We created the model and trained with 201 epochs. The latent representation generated by the scVI pipeline was stored into the original Seurat object and UMAP embeddings were calculated with the Seurat function `RunUMAP`, using the scVI reduction (1:10 dims).

For integration with the python package `SAMap` (77), gene-gene bipartite maps for each pair of species were built from reciprocal BLAST results (`map_genes.sh` script on <https://github.com/atarashansky/SAMap>). The files used for this step were: *Pleurodeles waltl*: translated proteome obtained by *in silico* translation (with `TransDecoder`) of the reference transcriptome described above, and annotated using `EGGNoG Mapper` as described above; turtle: RefSeq proteome of the turtle *Chrysemys picta bellii* downloaded from NCBI on May 17th 2022; lizard: RefSeq proteome of the lizard *Pogona vitticeps* downloaded from NCBI on May 17th 2022, genome assembly pv1.1.0. `SAMap` was run with the function `sm.run(neigh_from_keys =`

{'pw':True,'pv':True,'cp':True}), where the keys were clusters identities in salamander, lizard, and turtle.

For the cross-species label transfer results in Fig. S15, the label transfer algorithm in the R package Seurat was used, with the salamander neuronal data as "reference" and mouse data as "query". Two independent analyses were conducted, using mouse data from Zeisel et al. (78) and from Yao et al. (43), respectively. The Zeisel et al. dataset was filtered to keep only telencephalic neurons, according to the authors' cluster annotation; see above for a description of the Yao et al. dataset. Only one-to-one orthologs, identified as described above, were used for this analysis. Data were normalized using the SCTransform function in Seurat, and transfer anchors were identified with the function FindTransferAnchors (reduction="cca"). Label transfer was obtained with the function TransferData, which produces a matrix with prediction scores for the matching of each cell in the query to each cluster in the reference. These prediction scores are plotted in the UMAP space in Fig. S15.

In situ hybridization and immunohistochemistry

Animals were deeply anesthetized, with adults submerged in 0.2% with MS-222 and larvae submerged in 0.04% MS-222. All solutions during tissue preparation for *in situ* hybridization were prepared cold in RNase-free solutions prepared with DEPC-treated H₂O. Adult animals were transcardially perfused with 10 mL PBS, followed by 10 mL 4% PFA in PBS, then decapitated. Brains were dissected out, and postfixed overnight at 4°C in 4% PFA in PBS. Postfixation was stopped in PTW (PBS with 0.1% Tween). Coronal sections were prepared using a Leica VT1200S vibratome (70 µm thickness) with the brain embedded in 4% low-melting agarose. Larvae were fixed overnight at 4°C in 4% PFA in PBS, washed in PBS and cryoprotected in 30% sucrose-PBS. After embedding in Tissue-Tek OCT compound (Sakura), 12-20 µm coronal sections were cut on a cryostat and mounted on glass slides.

Immunohistochemistry

The sections were blocked in Blocking Buffer (2.5% BSA, 2.5% sheep serum, 50 mM glycine) in PBST (PBS with 0.2% Triton), then incubated with mouse anti-NeuN (1:500, Sigma-Aldrich MAB377), mouse anti-SATB2 (1:50, abcam ab51502) and/or rabbit anti-SOX2 (1:500, abcam ab97959) in primary Ab solution (10 mM glycine, 0.1% H₂O₂ in PBST) 1-3 nights at 4°C. Importantly, our scRNAseq analysis and *in situ* hybridization indicate that abcam ab51502 mouse anti-SATB2 recognizes salamander SATB1 instead (*Satb1* and *Satb2* are two recent paralogs). Specifically, *Satb1* *in situ* hybridization produced staining in the LPa and LPp, consistent with co-expression of *Satb1*, *Reln* and *Lhx2* in the lateral pallium. Additionally, *Satb2* *in situ* hybridization did not result in any staining, consistent with the fact that *Satb2* was not detected in scRNAseq data in pallial neurons. Therefore, in the main text, we will be referring to the SATB2 antibody as SATB1.

The tissue was washed 5 x 15 min in PBST at RT, then incubated in goat anti-mouse IgG, goat anti-rabbit IgG conjugated to Alexa 488, Alexa 594, or Alexa 647 (1:500, Invitrogen) with DAPI (1:5000) in PBST overnight at 4 °C or for 2 hrs at room temperature. The tissue was washed 5 x 15 min in PBST at RT. The tissue was washed once more in PBST and mounted in DAKO fluorescent mounting medium (Agilent Technologies). Images were acquired using a confocal microscope (Zeiss LSM800) and processed in FIJI.

In situ hybridization on sections

Vibratome sections were postfixed overnight at 4°C in 4% PFA in PBS, washed in PTW (PBS with 0.1% Tween-20), and the pia was peeled off manually with Dumont 5SF forceps (FST). Floating tissue sections were permeabilized for 8 min in 10 ug/mL ProK, washed 2 x 2 min in 2 mg/mL glycine, and rinsed in PTW for 5 min. To block nonspecific binding, the sections were acetylated by incubation for 5 min in 1% TEA/PTW, followed by 5 min in 1% TEA and 3 uL/mL acetic anhydride in PTW, and a 5 min wash in PTW. This was followed by postfixation for 20 min at RT, and 3 x 10 min washes in PTW. All stock solutions were prepared in DEPC treated or Molecular grade water.

To prepare probes for ISH, ~1kb fragments from coding regions of genes of interest were PCR amplified from a *Pleurodeles waltl* brain cDNA library or obtained by gene synthesis (Twist Bioscience) (see Data S1 for probe sequences). Fragments were cloned into the pCRII vector and sequences verified by Sanger sequencing (Eton Bio). After plasmid linearization, anti-sense DIG-labeled RNA probes were generated by *in vitro* transcription and purified with the RNeasy kit (Qiagen).

At 55-62°C, the tissue was pre-hybridized for 1 hr in Hybridization Mix (50% formamide, 5x SSC, 50 ug/mL heparin, 250 ug/mL yeast tRNA, 5x Denhardt's solution, 0.2% Tween, 500 ug/mL salmon sperm DNA, 10% Dextran sulfate in DEPC-treated water). The sections were then incubated 1-2 nights in Hybridization Mix with denatured riboprobes (1-3 ng/uL). Following hybridization, the sections were washed 2 x 30 min in a pre-warmed low stringency wash buffer (50% formamide, 2x SSC, 0.1% Tween in DEPC-treated water), 2 x 40 min in prewarmed high stringency wash buffer (0.1 - 0.2x SSC, 0.1% Tween), and 2 x 10 in MABT (100 mM maleic acid, 150 mM NaCl, pH 7.5) at RT. The sections were blocked in Blocking Buffer (5% sheep serum with 10% Roche Blocking Reagent 10x in MABT) for 1 hr at RT and incubated overnight with anti-Digoxigenin-AP (1:4000, Roche 11093274910) in Blocking Buffer.

Signal was developed by washing 4 x 30 min in MABT, followed by incubation in fresh staining solution (alkaline phosphatase buffer pH 9.5, 4.5 ug/mL NBT, 3.5 ug/mL BCIP, 5% polyvinyl alcohol) for 1-5 days. The staining reaction was stopped in PBS pH 7.4, and the sections were mounted in DAKO fluorescent mounting medium (Agilent Technologies). Images were acquired using an upright brightfield microscope (Leica DMR with Basler color camera, ACCU-Slide MS software). Background was subtracted evenly across the section using Photoshop CS6, and air bubbles accidentally present outside of the tissue section (for example, in the ventricle) were cropped out.

Hybridization Chain Reaction (HCR) in situ hybridization on sections

HCR-3.0-style probe pairs for fluorescent *in situ* mRNA detection were ordered from Molecular instruments or designed using the *insitu_probe_generator* (79) and ordered from IDT (20-33 pairs per probe set) (see Data S1 for probe sequences). After the same tissue pretreatment described above, the Molecular Instruments HCR v3.0 protocol for sample in solution (rev. 6) was followed (80). At 37°C, the tissue was pre-hybridized for 30 min in probe hybridization buffer (Molecular Instruments), and incubated overnight in probe hybridization buffer with 4- 20 pmol of probe sets (*Slc17a6*). Excess probe was removed by 4 x 15 min washes in probe wash buffer (Molecular Instruments) at 37°C, and 3 x 5 min washes in 5x SSCT at RT. The sections were preamplified for 30 min in amplification buffer (Molecular Instruments), and incubated overnight in amplification buffer with snap-cooled hairpins (60 pM, Molecular Instruments) in

the dark at RT. Excess hairpin was removed by 2 x 5 min, 2 x 30 min, 1 x 5 min washes in 5x SSCT at RT. Sections were incubated 1 hr in 1:5000 DAPI in 5x SSCT, washed 3 x 5 min in 5x SSCT, and mounted in DAKO fluorescent mounting medium (Agilent Technologies). Images were acquired using a confocal microscope (Zeiss LSM800) and processed in Fiji.

Whole-mount Hybridization Chain Reaction (HCR) in situ hybridization and iDISCO brain clearing

Tissue staining clearing was based on iDISCO (see (81)) and DIIFCO (HCR, see (80, 82)) protocols. Following fixation, the brains were rinsed 3 x 15 min in PBS. For clearing, the tissue was treated with an increasing gradient of methanol in PBS (20%, 40%, 60%, 80%, 100%, 1 hr each), bleached overnight with 5% hydrogen peroxide in 20% DMSO/methanol at 4°C, then rehydrated with a decreasing gradient of methanol. At 37°C, the brains were permeabilized for 1d in 0.2% Triton X-100, 20% DMSO and 0.3M glycine, the tissue was pre hybridized for 30 min in probe hybridization buffer (Molecular Instruments), and incubated 2 nights in probe hybridization buffer with 2-4 pmol of each probe set (*Etv1*, *Sox6*, *Slc17a6*, *Nr2f2*, *Penk*, *Rorb*). Excess probe was removed by 3 x 1 hr washes in probe wash buffer (Molecular Instruments) at 37°C, and washed overnight in 5x SSCT at RT. The tissue was preamplified for 30 min in amplification buffer (Molecular Instruments), and incubated 2 nights in amplification buffer with snap-cooled hairpins (60pM, Molecular Instruments) in the dark at RT. Excess hairpin was removed by 3 x 1 hr washes, followed by an overnight wash in 5x SSCT. The samples were embedded in 4% agarose, dehydrated in a decreasing methanol gradient, and incubated in 66% DCM/ 33% methanol for 3 hr at RT. Residual methanol was removed with 2 x 15 min washes in 100% DCM, and the tissue was allowed to clear overnight in DBE. Images were acquired using a LaVision Ultramicroscope II light sheet microscope at 4X magnification and 2 μm resolution, and subsequently visualized using ImarisViewer 9.8.0.

Axonal tracing

Injections and tissue processing

To provide access to all injection sites, and in accordance with similar tract-tracing experiments in amphibians (16), all injections were conducted *ex vivo*, under which conditions *Pleurodeles* brain tissue survives up to 2 days. Prior to *ex vivo* tracer injection, animals were deeply anesthetized by submersion in 0.2% MS-222. Animals were perfused transcardially with ice-cold oxygenated Amphibian Ringer's solution (96 mM NaCl; 20 mM NaHCO₃; 2 mM KCl; 10 mM HEPES; 11 mM glucose; 2 mM CaCl₂; 0.5 mM MgCl₂), and then decapitated. Brains were subsequently exposed while inside the skull, and the dura mater and choroid plexus removed. The arachnoid and pia mater were then removed using fine forceps only above the brain region to be injected. 10% 3 kD or 10 kD biotinylated dextran amines (3 kD BDA, Invitrogen D7135; 10 kD BDA Invitrogen D1956; diluted in 0.9% NaCl) were used as retrograde and anterograde tracers, respectively. While the two molecular weights do not result in exclusive unidirectional transport, the robust differential labeling of cell bodies (retrograde) versus axons (anterograde) depending on the injected molecule was clear during image analysis, and has been used historically (83, 84). All tracer solutions were pressure injected (20-30 nL, 1 nL/sec) using a glass capillary needle (diameter ~10 μm) connected to a Nanoject III injection system (Drummond); after injection, the needle was left in the injection site for 5 min to prevent leakage. Injection sites (OB (n=2), LPa (n=2), DP (n=2), MP (n=2), LA (n=2), VP_a (n=4), , VMH (n=1)) were determined using anatomical landmarks, and successful tracer application was

confirmed with FastGreen dye (Sigma-Aldrich F7252). Injected brains were transferred to fresh ice-cold Amphibian Ringer's solution and allowed to incubate for 24-48 hr on ice, with constant oxygenation. Brains were then fully extracted from the skull, transferred to 4% PFA, and fixed overnight at 4 °C. After fixation, 70 µm sections were cut on a vibratome, and processed either according to IHC or HCR staining protocol, as described above. During secondary antibody incubation (IHC) or during hairpin amplification (HCR), 1:500 Streptavidin conjugated to Alexa Fluor 488, 594, or 647 (Invitrogen S32354/S32356/S32354) was added to the solution to visualize BDA localization. All slices were subsequently mounted onto glass slides with DAKO fluorescent mounting medium (Agilent Technologies). Images of each section were acquired using a confocal microscope (Zeiss LSM800) and labeling was scored manually by three independent researchers. Localization and annotation of BDA tracer signal was determined using co-staining for known molecular markers, as well as anatomical landmarks (Fig. S1, Movie S1).

Supplementary text

Note on the nomenclature of telencephalic regions

The nomenclature of pallial regions in non-mammalian vertebrates has changed in the last decades as their identification and comparison improved over time. In amphibians, three pallial regions were initially described (medial, dorsal, and lateral pallium). However, neuroanatomical evidence for further subdivisions of the lateral pallium were also noticed (see for example (9, 22, 37)). With the introduction of the tetrapartite model for the pallium (17), the classical lateral pallium was subdivided in a lateral and a ventral pallium.

The major classes of pallial neurons identified in this work can be mapped to these four pallial regions (Fig. 2). Therefore, we are using the nomenclature that was established with the tetrapartite model for practical reasons, to keep references to these brain regions consistent with previous literature. However, the use of this terminology is not intended to endorse the tetrapartite pallium model, a model that aims to explain the evolutionary relationships of pallial regions in amniotes. The model has changed in recent years, and it implies evolutionary relationships that are not validated by our data. For example, in the recent version of the model, the term lateral pallium refers to the claustrum and insular cortex of mammals, and its putative sauropsid counterparts (85). In contrast, our data point to the homology of the amphibian lateral pallium (anatomical region) and olfactory-recipient cells in the reptilian lateral cortex.

Cross-species comparison of single-cell RNA sequencing data

Choice of algorithm

In order to compare cellular transcriptomes across species, we opted for manifold integration methods instead of the simpler comparison of pairwise cluster correlations across species (7). Integration methods have three advantages: (i) they exploit the high-dimensional gene correlations in scRNAseq datasets, instead of “flattening” the comparisons to pairwise correlation scores; (ii) they preserve the single-cell structure of the data, instead of averaging information across clusters; therefore, they are not sensitive to the depth of clustering (unless clustering results are used in downstream analysis); (iii) they treat the “species signal” as a “batch effect”, removing all gene expression variation that arises from concerted gene expression evolution at the global scale of the organism (details below, (86)).

Data integration algorithms are typically developed and optimized for the comparison of samples from the same species (“treated” vs “control”) or closely-related species (“mouse” vs “human”). The integration of data across distantly-related species poses some specific challenges:

- Gene choice: the most common choice to select genes for alignment is to choose variable genes among one-to-one orthologs (Seurat, Harmony, scVI). One-to-many and many-to-many orthologs are excluded because duplicated genes may diverge in function, making their comparison problematic. The number of one-to-one orthologs correlates with the evolutionary distance of the species being compared, and therefore, the comparison of distantly-related species removes many species-specific paralogs that may play a critical role in cellular diversity.

In contrast, SAMap takes into account paralogs, motivated by the observation that paralogs may replace each other’s expression in homologous cell types and tissues as species diverge (“paralog switching”). To include paralogs for scRNAseq data integration, SAMap weighs the edges of the nearest-neighbor graph using reciprocal blast scores as a proxy for gene similarity (77).

- “species signal”: transcriptomics data, regardless of the tissues or cell types compared, cluster by species, unless a batch correction algorithm is applied. The so-called species signal (86) arises from global differences of gene expression, caused by drift or by pleiotropic changes of gene regulatory networks. This component of gene expression variation is not associated with selective processes that act on specific cell types or tissues under comparison. Correcting for this component of gene expression variation ensures accurate cross-species comparisons. However, single-cell integration algorithms treat and overcome the batch correction problem in different ways, and therefore, their ability to correct the species’ signal is expected to differ across methods (87).

Results and their interpretation

We used Seurat (CCA), Harmony, SAMap, and scVI to perform the integration of salamander, lizard and turtle telencephalic data (Fig. 4 and Fig. S11) and of salamander, lizard, turtle and mouse medial and dorsal pallium cells (Fig. 5 and Fig. S17). These analyses were performed following the recommendations of the developers. Therefore, data normalization, choice of variable genes, and other parameters were largely independent between the analyses (see Methods for details). Results were plotted as UMAP embeddings. Seurat, Harmony, and SAMap all produced integrations with good species mixing. In contrast, species segregation was more pronounced in scVI embeddings, suggesting that this algorithm is more conservative and does not correct entirely for the species signal. Coloring the UMAPs by telencephalic region showed that the four algorithms group together cells sampled from homologous brain regions in a way that is consistent across algorithms.

It should be stressed that these integration algorithms are designed to capture common gene expression variation, and they need to be analyzed carefully before drawing evolutionary inferences from them:

- *co-clustering of cells from different species does not imply that cells have identical transcriptomes.* The alignment is based on genes that have shared patterns of co-variation across the datasets, therefore it brings together cells from different species as long as a substantial fraction of their transcriptome is similar (after batch correction). Hundreds of genes are differentially expressed in cells from different species that co-cluster, even in human-mouse comparisons (88).

- *co-clustering of cells from different species does not imply cell type homology*. Gene expression similarity can arise from homology or gene expression convergence; the latter case is when the same effector genes are expressed under the control of different sets of transcription factors. For this reason, transcriptomic similarities cannot be taken as the exclusive evidence for homology. The comparison of transcription factor expression, and information on the topological position of neurons in the brain and on neuronal connectivity should be used to support or refute homology hypotheses. For example, we interpret the co-clustering of the salamander POE and reptilian rostral aDVR as a case of gene expression convergence, because these two structures (i) express different sets of transcription factors (ii) have different topological positions in the pallium (POE is anterior and medial, aDVR is more posterior and lateral) and (iii) have different patterns of connectivity.

Supplementary figures

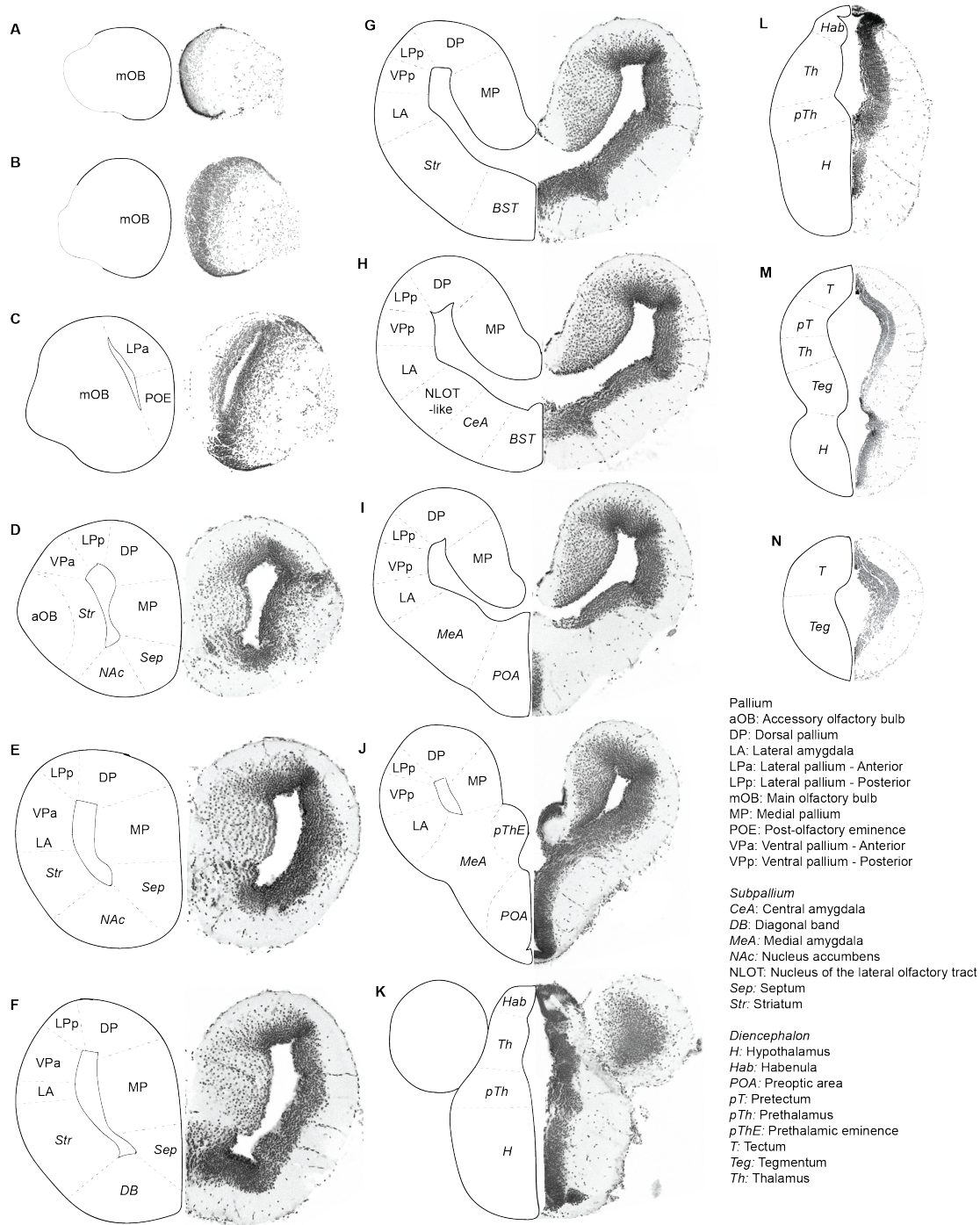


Fig. S1. DAPI atlas of the *Pleurodeles* brain. Coronal sections (70 μ m thick) of the *Pleurodeles waltl* brain, arranged anterior (A) to posterior (N). Nuclear (DAPI) stain is shown on the right in each panel, with a schematic showing subdivisions on the left, annotated with the acronyms described in the key. Annotations were performed using histological data from the present study, in tandem with previous literature (9, 10, 89–92). See also Movie S1.

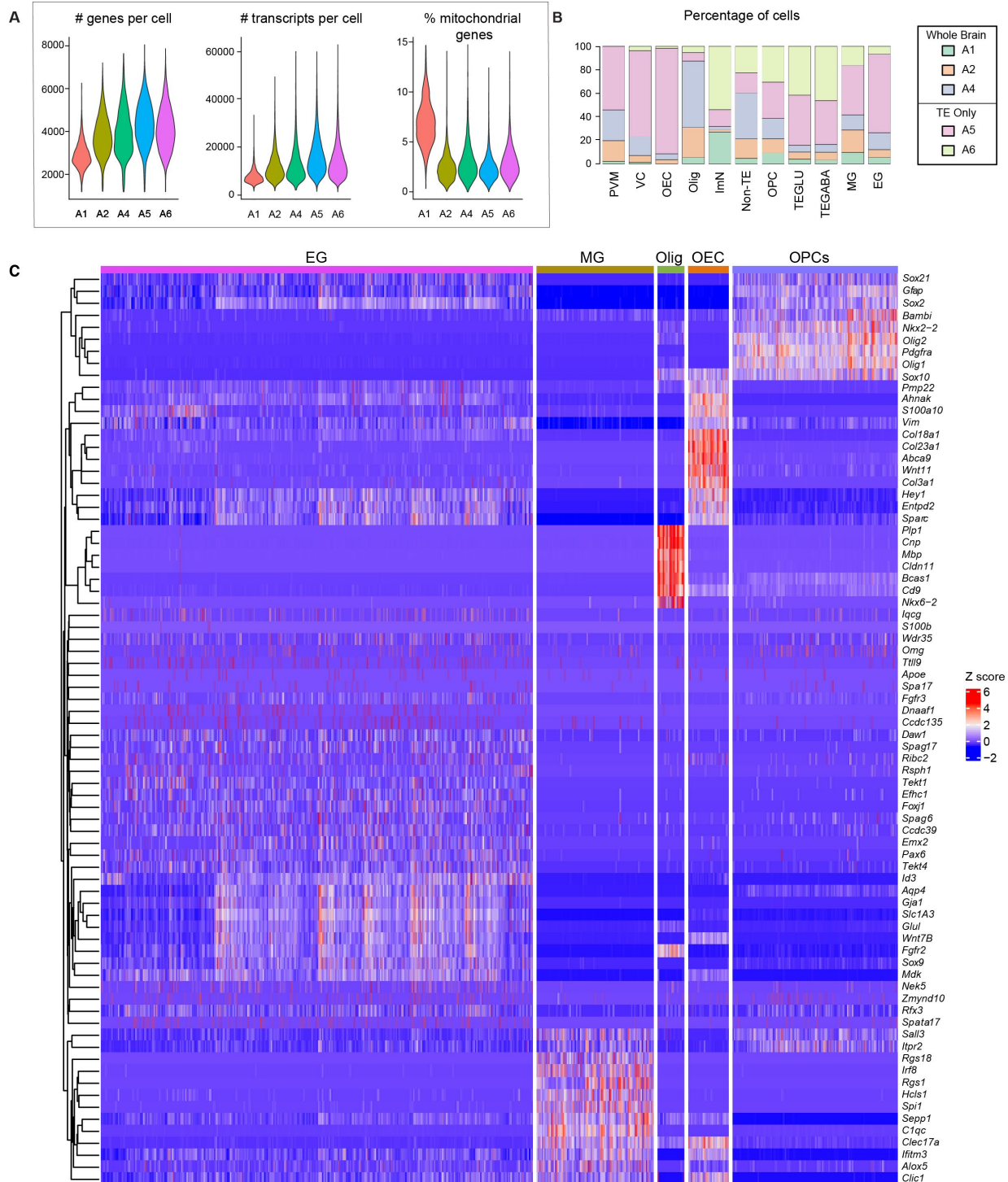


Fig. S2. Quality control of the adult brain single-cell dataset and additional gene expression profiles of non-neuronal cells. (A) Violin plots of the number of genes or transcripts per cell, and percent mitochondrial genes in the global dataset, sorted by animal. **(B)** Bar plot of percentage of cells per animal in global dataset clusters. Legend denoting tissue included for each animal in single cell prep for dissociation. TE: telencephalon. **(C)** Left: dendrogram hierarchically clustered based on molecular similarity. Right: heatmap of differentially expressed genes for non-neuronal clusters.

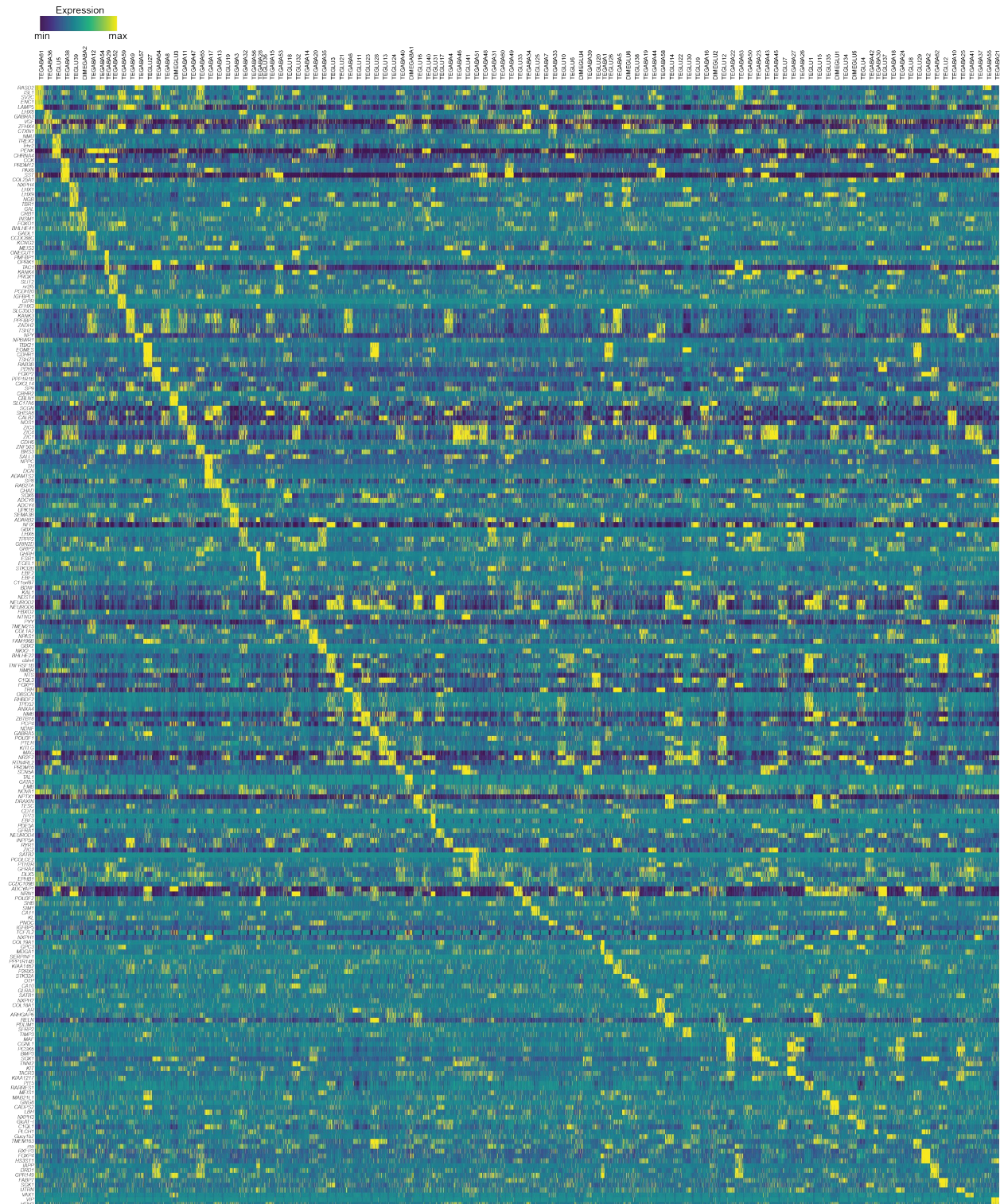


Fig S3. Cluster marker genes in the *Pleurodeles* neuronal dataset. Heatmap showing the expression (Z score) of the top 5 cluster-specific marker genes in the salamander neuronal dataset (rows: genes, columns: single cells grouped by cluster membership, up to 50 cells for each cluster, randomly selected).

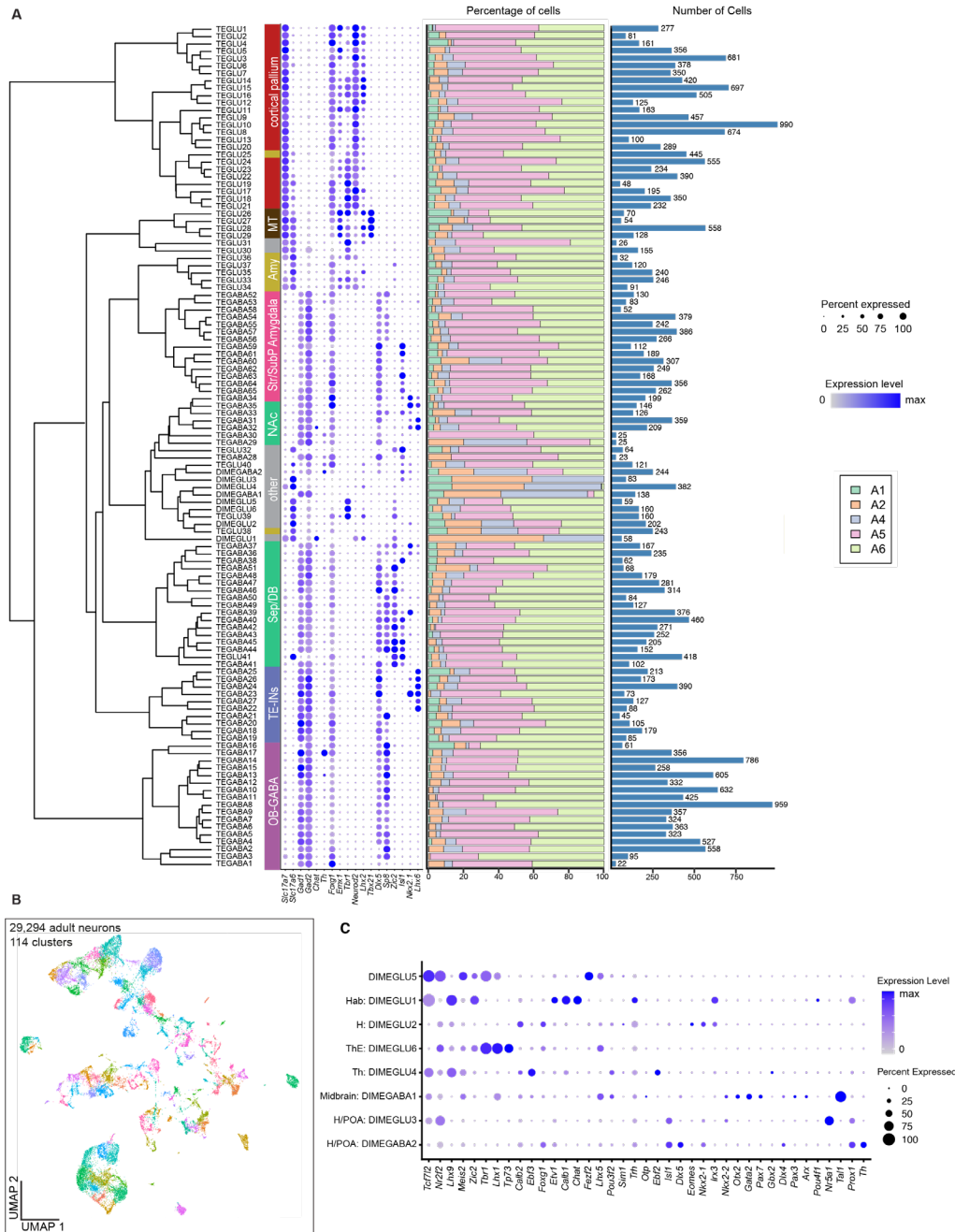


Fig. S4. Cluster annotation of the *Pleurodeles* neuronal dataset. (A) From left to right: hierarchical clustering of cluster average gene expression profiles. The dendrogram was built on the basis of the top 2000 variable genes, and groups together clusters with similar gene expression profiles. These groups correspond largely, but not exactly, to anatomically-defined regions of the salamander telencephalon. To annotate the regional identity of these clusters and groups, we analyzed the expression of regional markers, such as those shown in the DotPlot in the middle. Right: Bar plot of percentage of cells per animal in each cluster; Bar plot of total number of cells in each cluster. **(B)** UMAP representation of 29,294 adult salamander neuron transcriptomes, colored by cluster identity. **(C)** DotPlot of marker genes (columns) used for the annotation of neurons in the di- and mesencephalon ("DIME", rows). See Figure S1 for abbreviations.

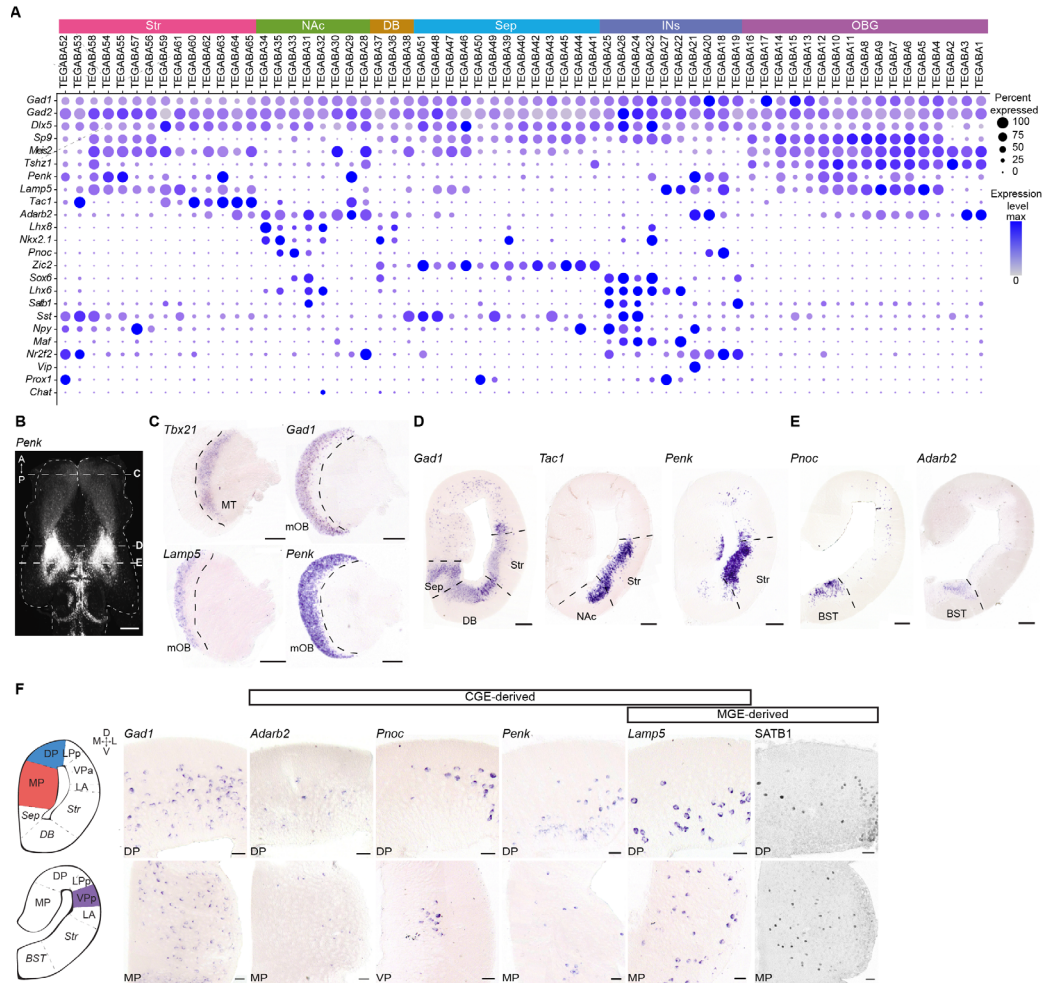


Fig. S5. GABAergic neuron types in the *Pleurodeles* telencephalon. (A) DotPlot showing the expression of key marker genes defining subpallial regions and telencephalic interneurons. (B) Whole mount HCR-ISH for *Penk* with dashed lines indicating the sectioning plane for C-E. Scale bar: 500 μ m. (C) Coronal sections through the olfactory bulb showing mitral and tufted cells (*Tbx21*+) and olfactory bulb interneurons (*Gad1*+) . Subpopulations of olfactory bulb interneurons express *Lamp5* and/or *Penk*. (D) Coronal sections through the telencephalon showing the expression of *Gad1* (all GABAergic neurons), *Tac1* (striatum, accumbens), and *Penk* (striatum and subset of pallial interneurons). (E) Coronal sections through the telencephalon showing the expression of *Pnoc* and *Adarb2* (BST and pallial interneurons). (F) Magnifications from coronal sections through the telencephalon showing sparse MGE- and CGE-derived interneurons spread throughout the pallium. Schematics on the left show magnified pallial regions. Marker genes were selected from the scRNAseq data for their specific expression in interneuron types. MGE- and CGE-derived interneurons can be found scattered throughout the pallium, including the dorsal and the medial pallium as shown here. Some interneuron subtypes, such as the CGE-derived *Adarb2*+ *Penk*+ neurons, were enriched in DP in comparison to MP. See methods for specifics on SATB1 antibody. Abbreviations: BST, bed nucleus of the stria terminalis; CGE, central ganglionic eminence; DB, diagonal band; DP, dorsal pallium; INs, interneurons; MGE, medial ganglionic eminence; MP, medial pallium; NAc, nucleus accumbens; OBG, olfactory bulb GABAergic; Sep, septum; Str, striatum; VP, ventral pallium.

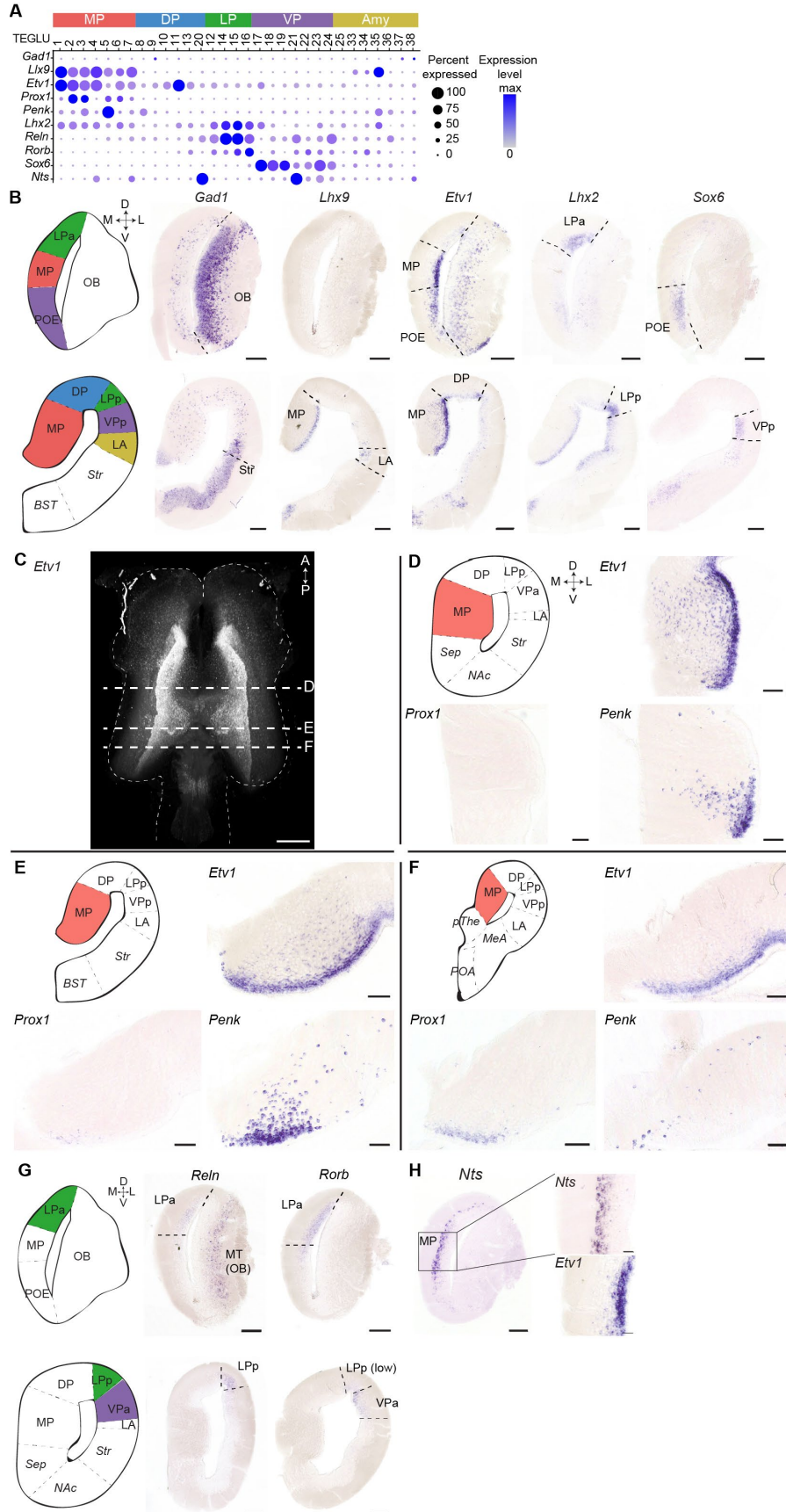


Fig. S6. Regionalization of the pallium in *Pleurodeles*. (A) DotPlot showing the expression of the genes shown in panels B-H. (B) Left to right: schematics of coronal sections through the telencephalon at different anterior-posterior levels; expression of *Gad1*, indicating the pallium-amygdala-subpallium boundaries, and of transcription factors labeling distinct pallial regions along the mediolateral axis: expression of *Lhx9* in the MP and in cells of the lateral amygdala; expression of *Etv1* showing the transition between medial and dorsal pallium (TEGLU 8-11, 13, 20), clearly demarcated by a ventricular sulcus and an abrupt change of cell density; expression of *Lhx2* showing the lateral pallium which is bordered by the axons of the lateral olfactory tract; expression of *Sox6* in the ventral pallium which is molecularly diverse, in line with its anatomical heterogeneity. Scale bars: 200 μ m. (C) Whole mount HCR-ISH for *Etv1* with dashed lines indicating the sectioning plane for D-F. Scale bar: 500 μ m. (D-F) Colorimetric ISH on floating tissue sections showing expression of *Etv1*, *Prox1* and *Penk* across the anterior-posterior axis in the medial pallium. *Prox1*, a marker of the dentate gyrus (DG), overlaps with *Etv1*, a marker of the Cornu Ammonis (CA) field, therefore showing no discrete subdivisions resembling CA or DG subfields in salamander. Scale bars: 50 μ m. (G) Left: schematics of coronal sections at different anterior-posterior levels. Right: Expression of *Reln* in the olfactory bulb mitral and tufted cells, and in the anterior and posterior lateral pallium (aLP and pLP), and of *Rorb* in the anterior and posterior LP, and in the anterior VP. Scale bars: 200 μ m. (H) Expression of *Nts* and *Etv1* in the anterior medial pallium showing the presence of at least two molecularly different layers in the pallium. Scale bars in right panels: 50 μ m. For abbreviations, see Fig. S1; A, anterior; D, dorsal; L, lateral; M, medial; P, posterior; V, ventral.

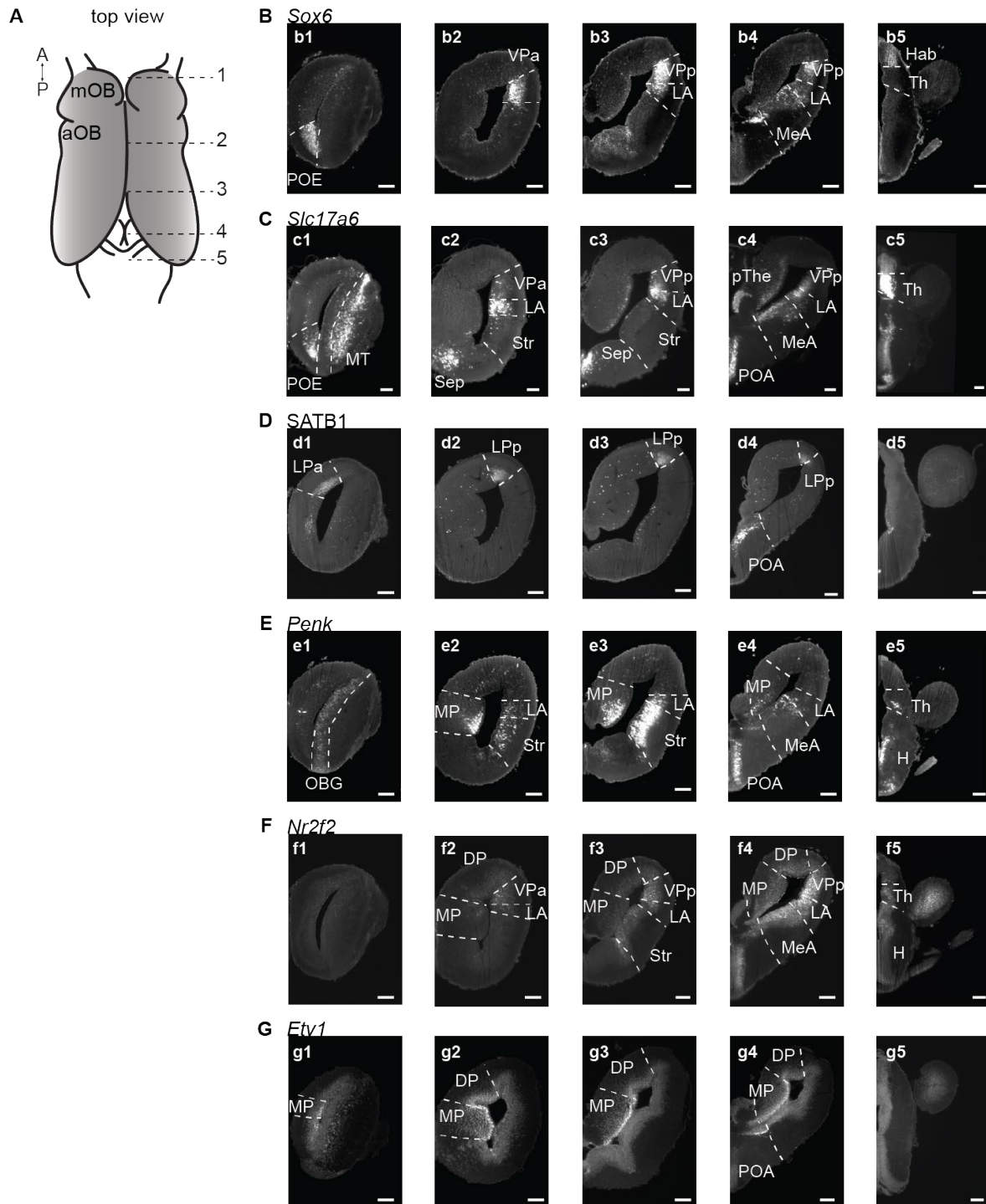


Fig S7. Expression of marker genes detected by whole-brain HCR *in situ* hybridization in the *Pleurodeles* forebrain. (A) Top view schematic diagram of the *Pleurodeles* forebrain, dotted lines represent optical planes of respective numbered images in panels B-G. (B-G) Optical coronal sections from whole cleared brains imaged by light sheet microscopy, following HCR *in situ* hybridization (B, C, E, F, G) or immunostaining (D). Dotted lines represent boundaries of telencephalic or diencephalic territories. See Movies S2-3 for the corresponding volumetric data. For abbreviations, see Fig S1. Scale bars: 200 μ m.

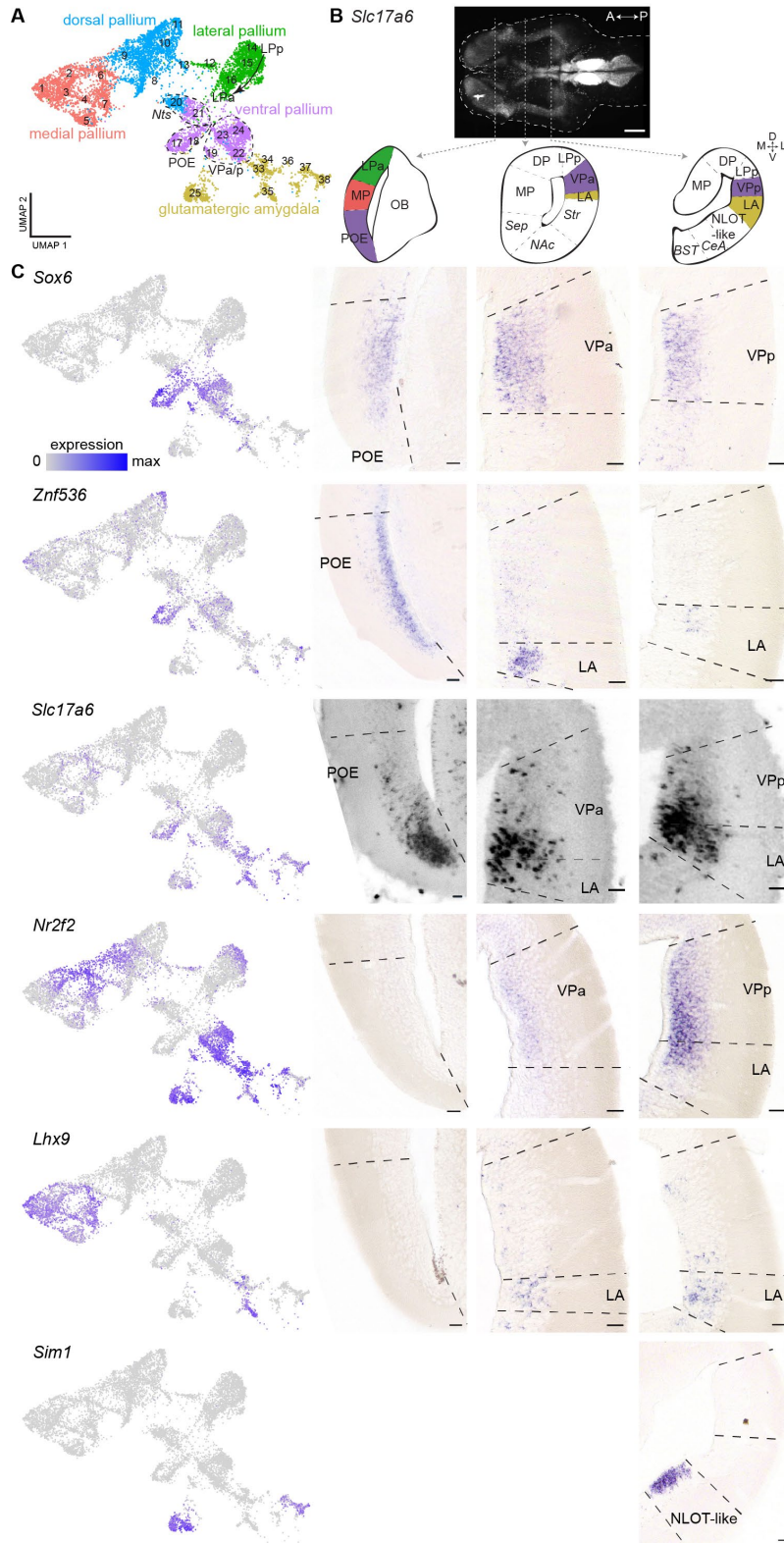


Fig. S8. Regionalization of the ventral pallidum and pallial amygdala in *Pleurodeles*. (A) UMAP representation of cortical pallium and amygdala neurons, color coded according to brain region. (B) Whole mount HCR-ISH for *Slc17a6* and schematic overviews of the telencephalon at

three levels across the anterior-posterior axis, with indication of the ventral pallium and lateral amygdala. (C) Left: UMAP plots of pallial glutamatergic neurons, with single-cells color-coded by the expression of marker genes *Sox6*, *Znf536*, *Slc17a6*, *Nr2f2*, *Lhx9* and *Sim1*. Right: gene expression on coronal sections through the telencephalon at three representative levels across the anterior-posterior axis. These marker genes define multiple subdivisions in the ventral pallium and amygdala, with an anterior *Znf536*-positive post-olfactory eminence, which is further subdivided in a more dorsal *Slc17a6*-negative and a ventral *Slc17a6*-positive domain. At intermediate levels, *Nr2f2* expression increases and *Slc17a6* is expressed in a gradient from low dorsally to high ventrally and highest in the LA. At posterior levels, ventral from the *Sox6*-positive/*Nr2f2*-positive VPP, the LA is defined by high-level expression of *Slc17a6* and *Nr2f2*, expression of *Lhx9* in a subset of cells and absence of *Sox6*. Lastly, *Sim1*-positive cells demarcate the NLOT-like cells in the amygdala (see also (23) for a detailed description of amygdala neuron types in *Pleurodeles*). Scale bars: 50um. For abbreviations, see Fig. S1; A, anterior; D, dorsal; L, lateral; M, medial; P, posterior; V, ventral.

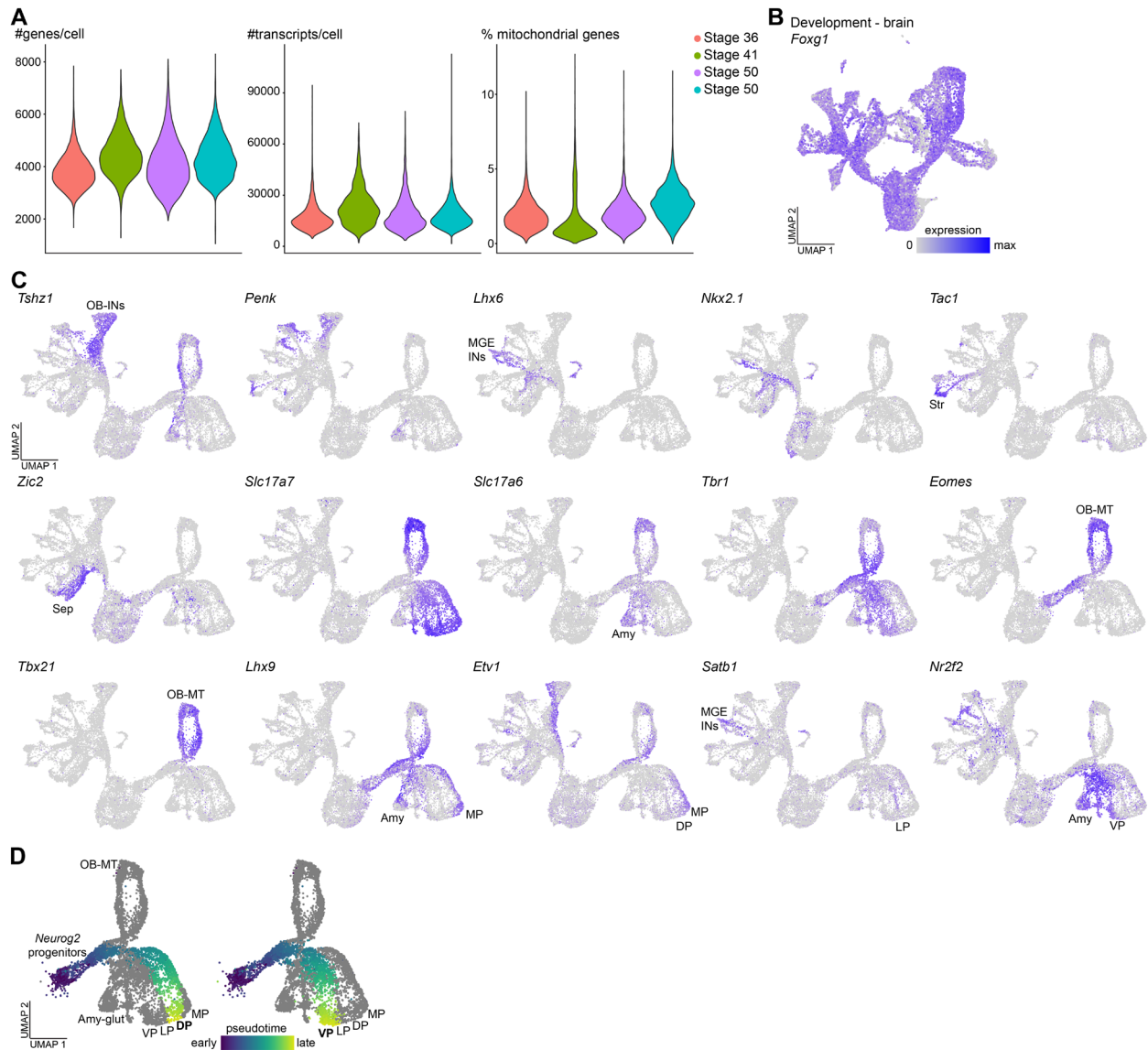


Fig. S9. Quality control and additional data for the developmental dataset. (A) Violin plots of the number of genes or transcripts per cell, and percent mitochondrial genes in the full developmental neuronal dataset prior to filtering for telencephalic clusters, violin plots sorted per stage. (B) Expression of *Foxg1* in the UMAP space, showing that some clusters express no or low levels of *Foxg1* and are therefore non-telencephalic. These clusters were removed from further analysis. (C) UMAP plots showing telencephalic cells colored by the expression of marker genes that were identified for their specific expression in telencephalic regions in the adult dataset. These data support the label transfer results in Fig. 3C. (D) Pseudotime values for DP and VP that were used to order cells for heatmap in Fig. 3F. See Methods for details on the pseudotime analysis. Abbreviations: Amy, amygdala; DP, dorsal pallium; INs, interneurons; LP, lateral pallium; MP, medial pallium; MGE, medial ganglionic eminence; MT, mitral and tufted cells of the olfactory bulb; OB, olfactory bulb; OBG, olfactory bulb GABAergic; Sep, septum; Str, striatum; VP, ventral pallium.

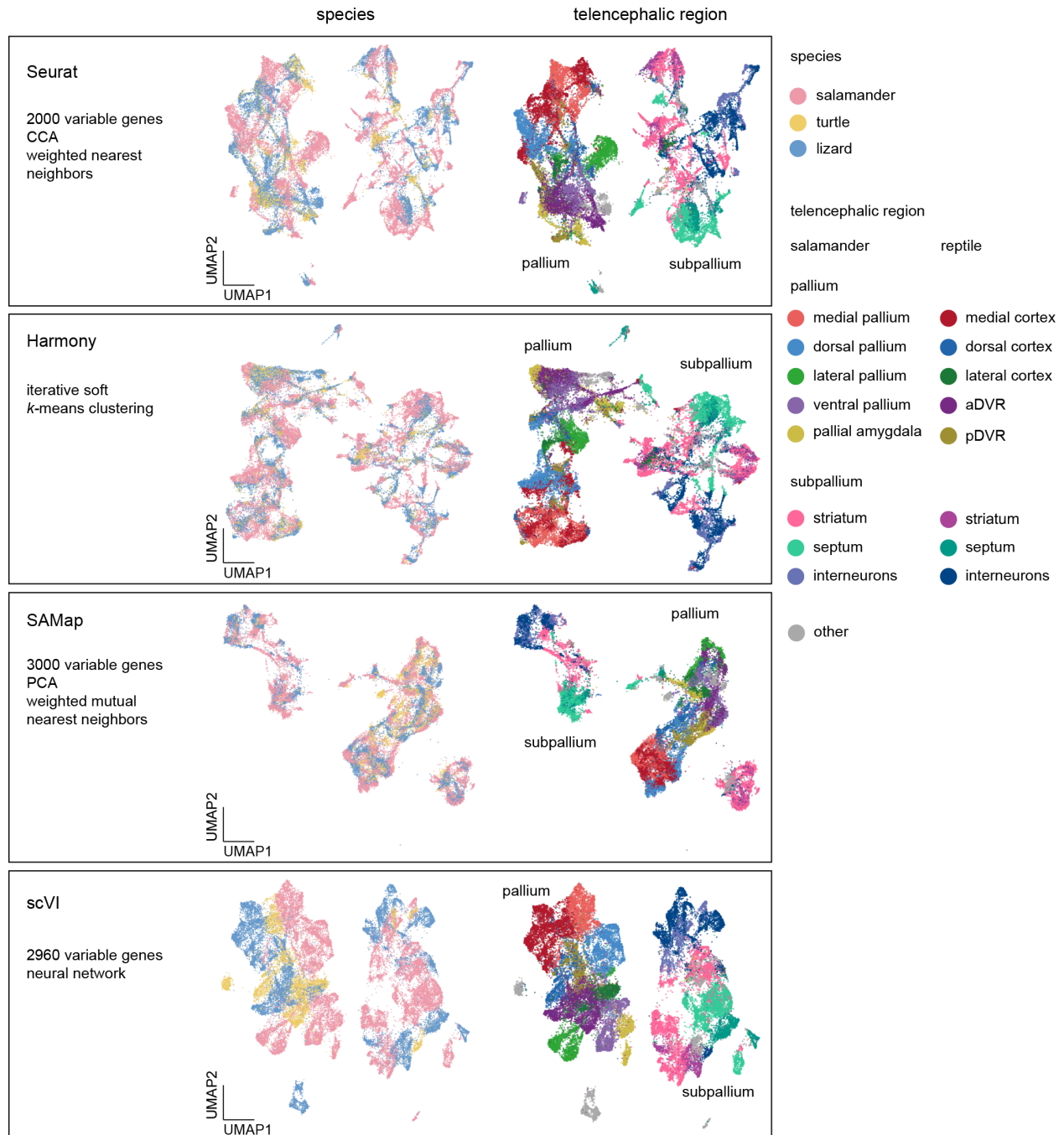


Fig. S11. Manifold integration of scRNAseq data from the telencephali of salamander, turtle, and lizard using four different integration algorithms (compare to Figure 4). From top to bottom: integration of the same datasets with Seurat (31), Harmony (75), SAMap (77), and scVI (76). Analysis parameters are indicated, more details available in Methods. First column: UMAP embeddings with cells color-coded by species. Second column: UMAP embeddings with cells color-coded by telencephalic region, in two different shades (salamander and reptile). Similar species mixing is obtained with Seurat, Harmony, and SAMap. Results from the Seurat integration were largely recapitulated by using alternative integration algorithms as the main telencephalic regions from the different species cluster closely together.

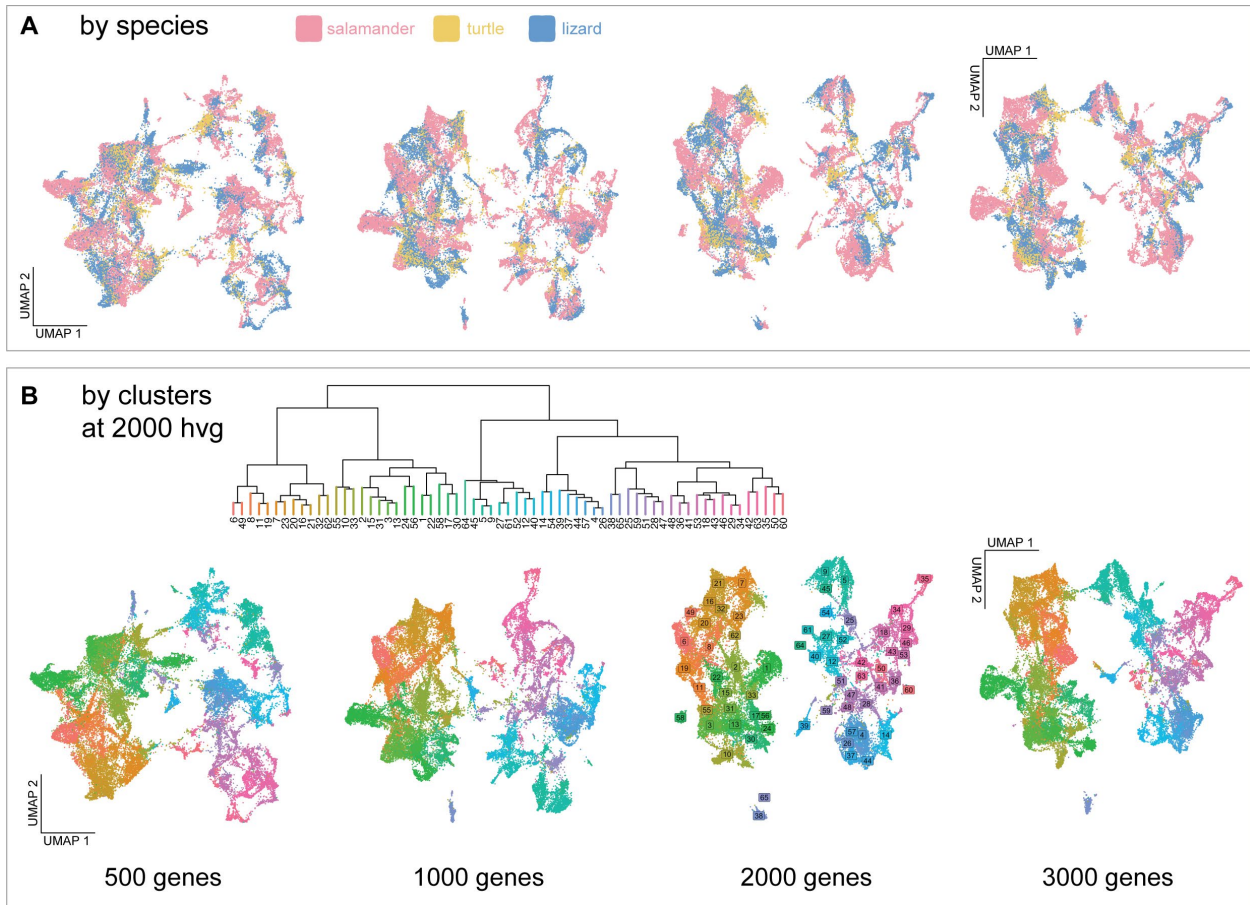


Fig. S12. Integration of salamander, lizard, and turtle data, using the Seurat pipeline with different numbers of high variable genes (hvg) (compare to Figure 4) (A) UMAP plots showing single cells colored by species after Seurat integration analyses run with 500, 1000, 2000, and 3000 hvg. Changing the number of hvg did not affect species mixing in the UMAP embedding. (B) Top: average gene expression profiles were computed for the clusters identified in the integrated space (2000 hvg), and then used for hierarchical clustering. The terminal branches of the resulting dendrogram were color-coded according to the position in the tree. Bottom: UMAP plots of Seurat integrations. Colors in the UMAP computed with 2000 hvg correspond to the color in the dendrogram. Single cells in the other UMAPs (500, 1000 and 3000 hvg) are colored with the same colors used in the reference integration (2000 hvg). This indicates that the high-level structure of the data manifold is relatively stable when different numbers of hvg are used.

Fig. S13. Integrated clusters of salamander, lizard, and turtle telencephalic neurons. Top: dendrogram of cluster average gene expression profiles. The dendrogram was built on the basis of the top 2000 variable genes, and groups together clusters with similar gene expression profiles with branches colored by species mixture (gray, equal proportion of cells from each species). Bottom: the heatmap under the tree shows the percentage of cells from each original cell cluster according to its species-specific annotation (rows) that land in the integrated clusters (columns). Our integrated dataset splits broadly into two groups defined by GABAergic and glutamatergic identity.

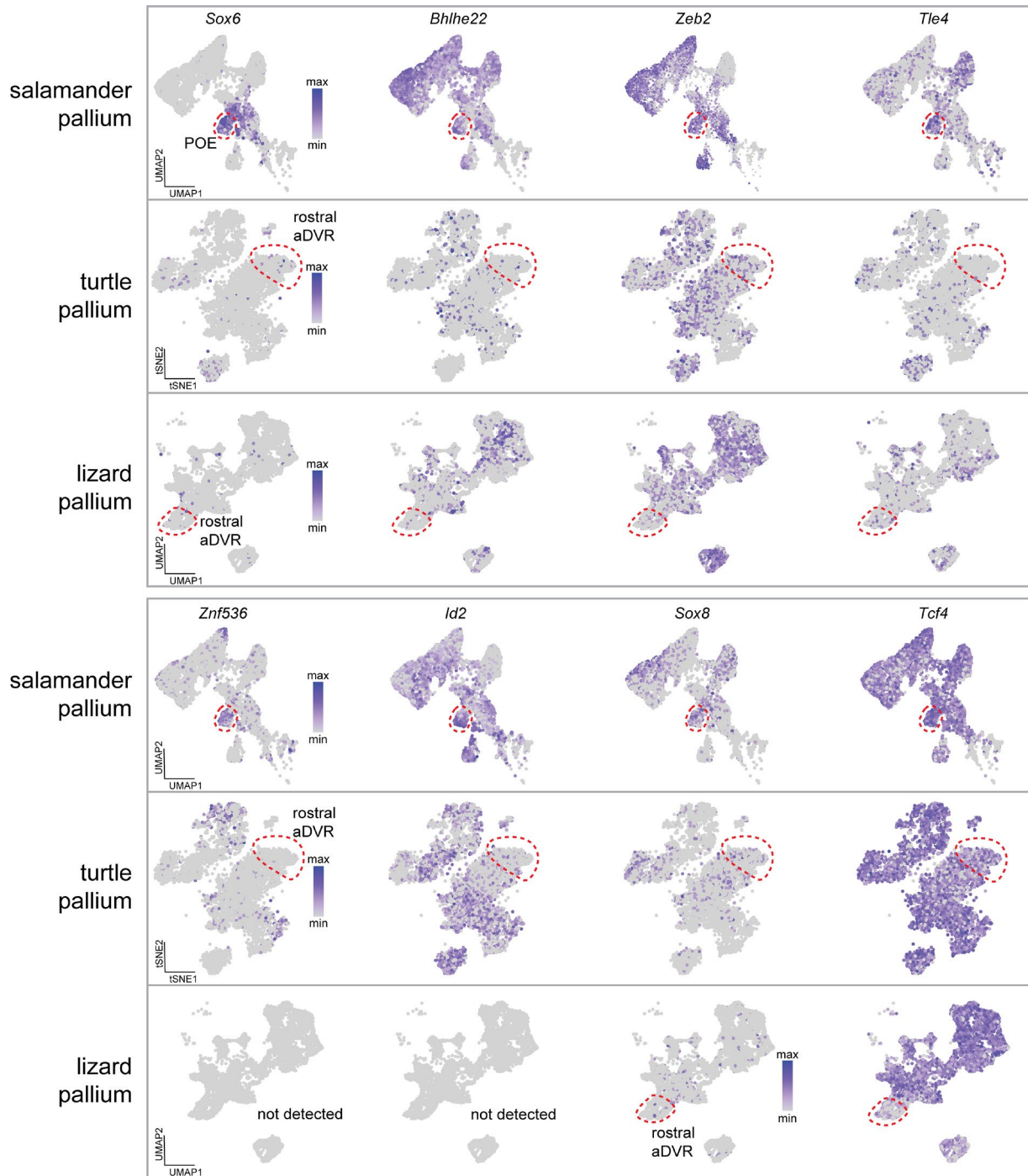


Fig. S14. Feature Plots of salamander POE and reptile aDVR transcription factors. UMAP plots showing salamander, lizard and turtle subsetted pallium datasets colored by expression levels of transcription factors that were identified for their differential upregulation in the salamander POE (circled in red, top panel). Expression of these transcription factors is low or not detected in turtle and lizard rostral aDVR (circled in red, middle and lower panels). This suggests that the clustering of POE/rostral aDVR in the integrated dataset is driven by effector genes and, thus, the rostral sensory aDVR evolved by recruiting effector genes used in other sensory-processing areas.

label transfer: salamander reference and mouse query

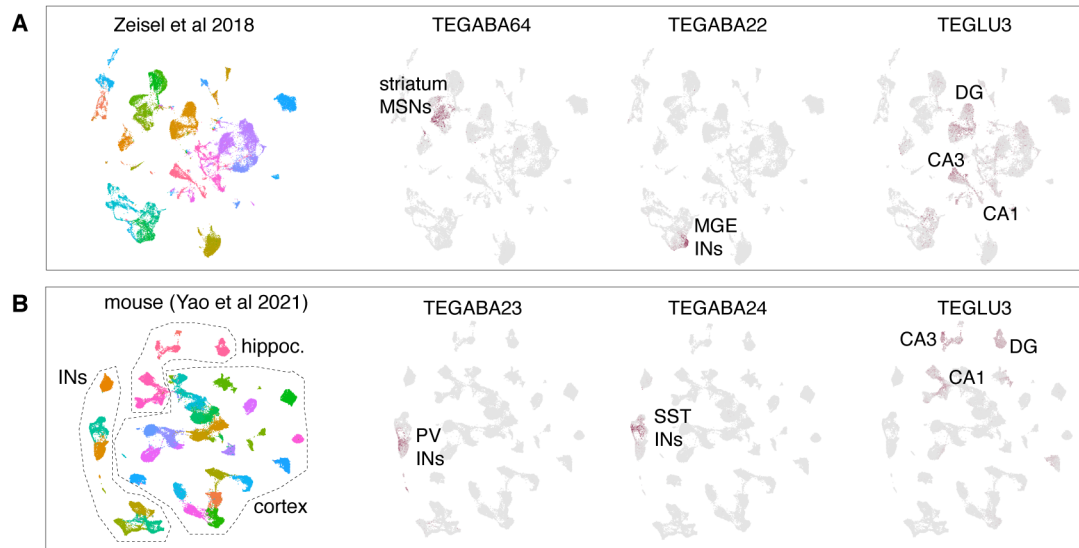


Fig. S15. Comparison of salamander and mouse scRNAseq data by label transfer. (A) UMAP plots of telencephalic neurons from Zeisel et al. 2018. Left: cells color-coded by cluster. Right: cells color-coded by the label transfer score for the salamander clusters indicated. Medium spiny neurons (MSN) in the mouse striatum mapped with high scores (label transfer score) on the salamander TEGABA64 cluster (striatum), mouse MGE interneurons mapped on salamander TEGABA22 cells, and mouse DG, CA1 and CA3 neurons mapped on salamander TEGLU3 cells. **(B)** Same as A, but using a downsampled Yao et al. 2021 dataset (cortex and hippocampus) as query. Mouse Pvalb (PV) interneurons mapped on the salamander TEGABA23 cluster, mouse SST interneurons on salamander TEGABA24, and mouse DG, CA1 and CA3 neurons on salamander TEGLU3.

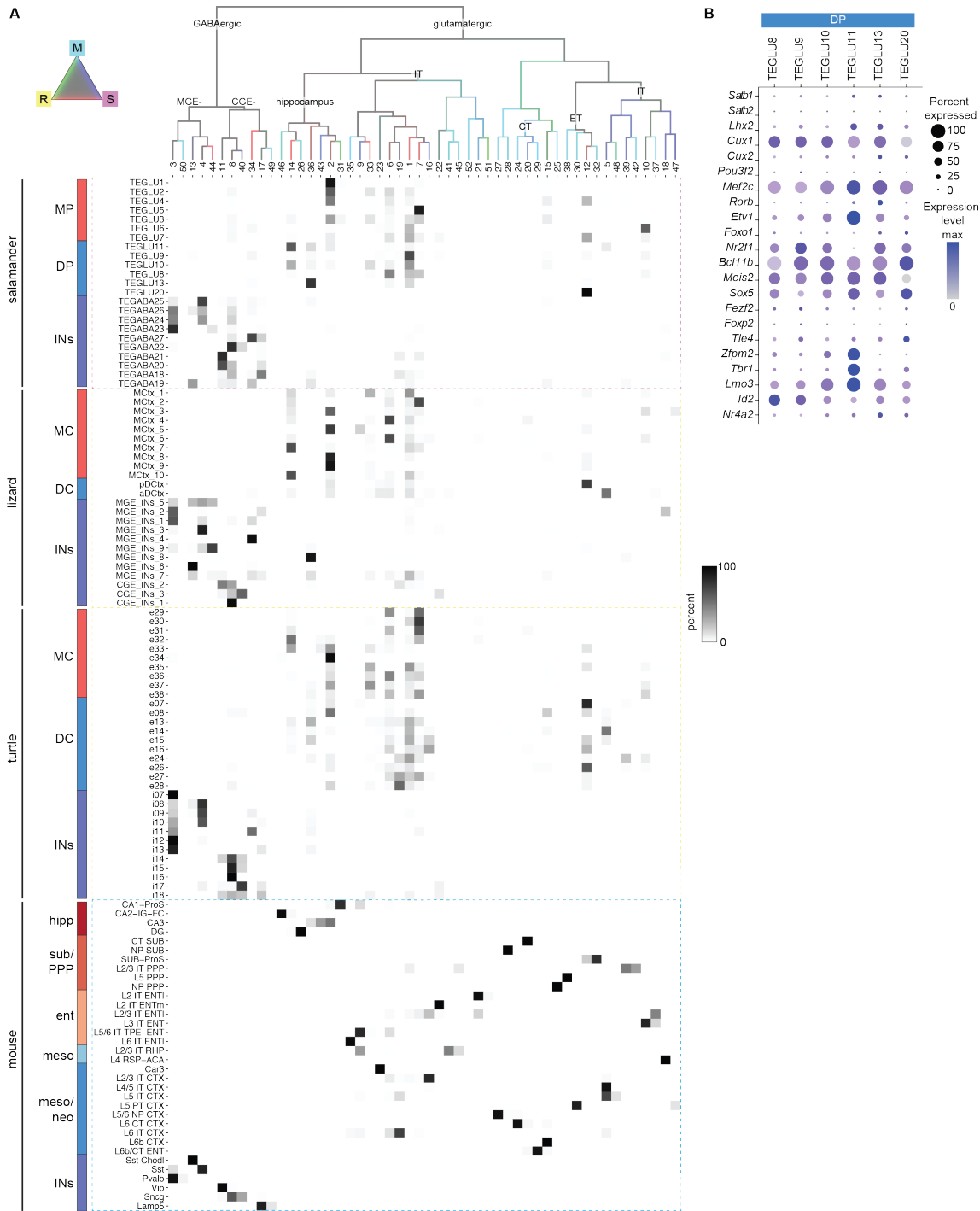


Fig. S16. Integration of salamander, turtle, lizard, and mouse neurons from the dorsomedial pallium. (A) Top: dendrogram shows the molecular similarity of integrated clusters, with branches colored by species mixture (gray, equal proportion of cells from each species). Bottom: the heatmap under the tree shows the percentage of cells from each original cell cluster according to its species-specific annotation (rows) in the integrated clusters (columns). **(B)** DotPlot showing the (absence of) expression of canonical mammalian layer transcription factors in salamander dorsal pallium clusters. Abbreviations: CT, corticothalamic; ET, extra telencephalic; IT, intratelencephalic.

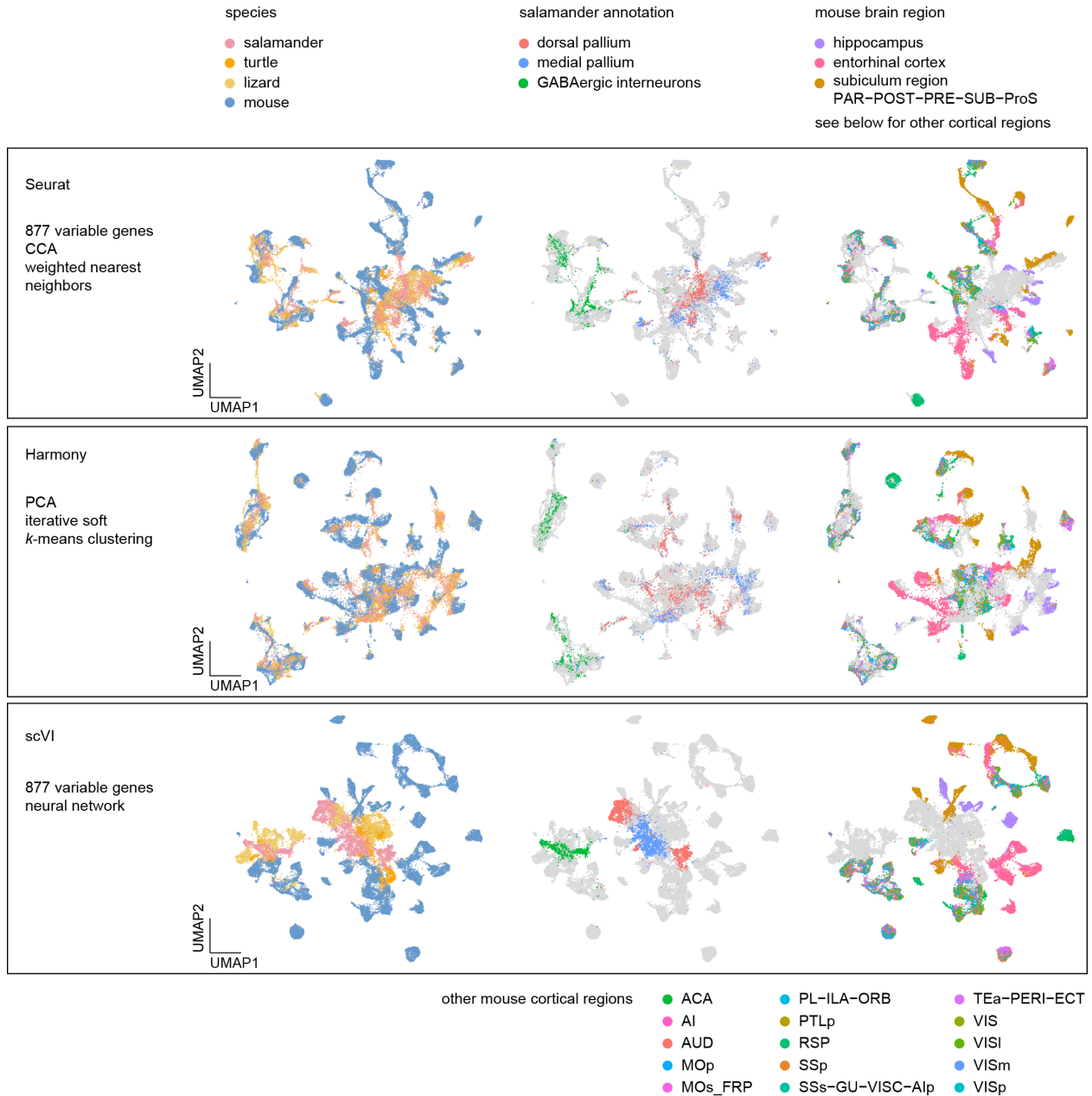


Fig. S17. Manifold integration of scRNAseq data from the telencephali of salamander, turtle, lizard and mouse using three different integration algorithms (compare to Figure 5). From top to bottom: integration of the same datasets with Seurat (31), Harmony (75), and scVI (76). Analysis parameters are indicated, more details available in Methods. First column: UMAP embeddings with cells color-coded by species. Second column: UMAP embeddings with cells color-coded by telencephalic region and cortical interneuron identities (salamander and reptile). Third column: UMAP embeddings with cells color-coded by mouse brain regions. Similar species mixing is obtained with Seurat, Harmony, and scVI. Results from the Seurat integration were largely recapitulated by using alternative integration algorithms as the main telencephalic regions from the different species cluster in a similar way.

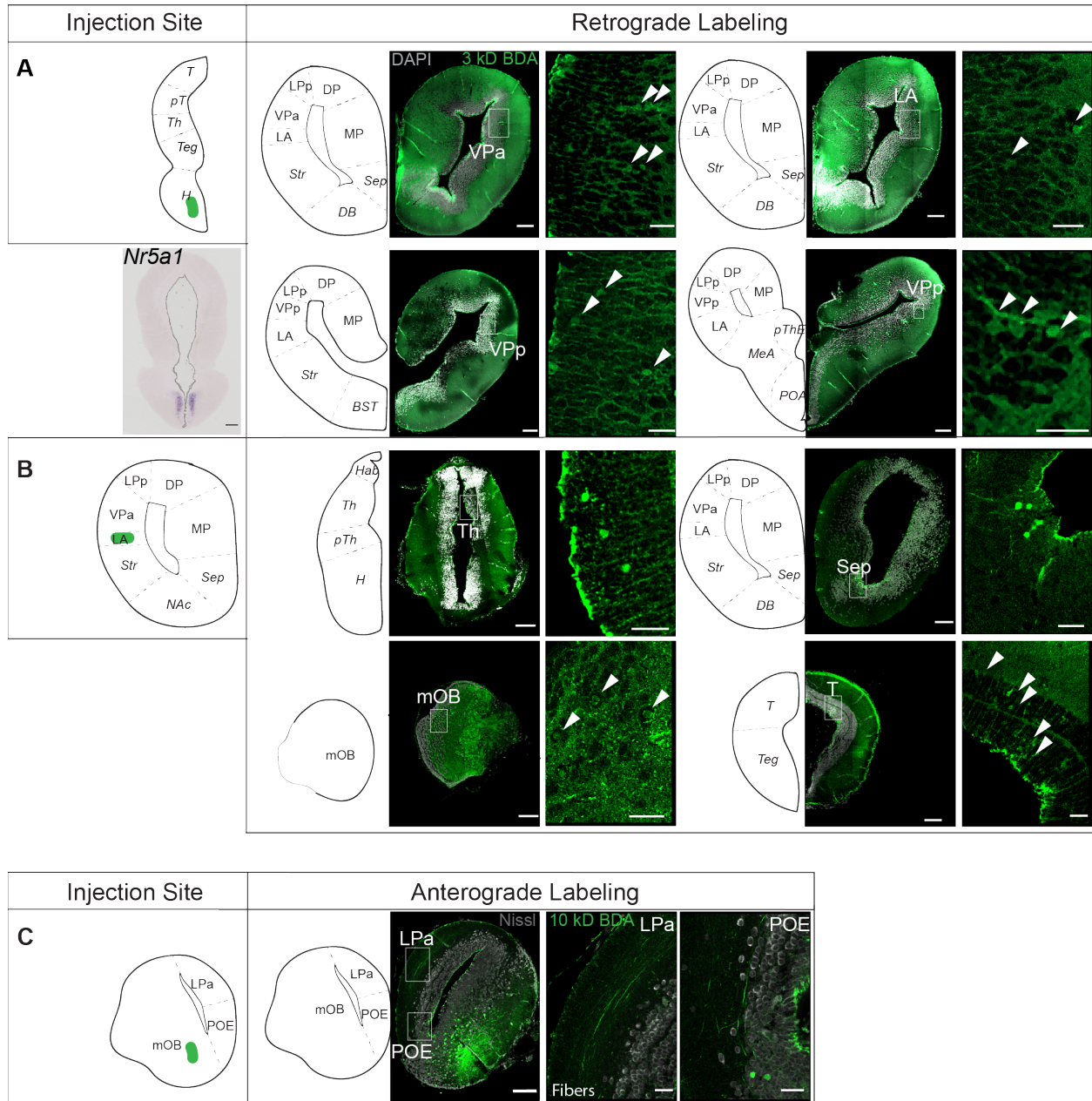


Fig. S18. Additional tracing data: injections of tracers in the hypothalamus, lateral amygdala, and olfactory bulb. Representative images of BDA tracer injections and labeling of (A) VMH injection (n=1) and retrogradely labeled cells with *Nr5a1* *in situ* hybridization; (B) lateral amygdala injection (n=2) and retrogradely labeled cells; (C) mitral tufted cells of the mOB and anterogradely labeled fibers projecting to the lateral pallium and to the postolfactory eminence (n=2). Scale bars: 200 um for slice overview, 50 um for magnified images.

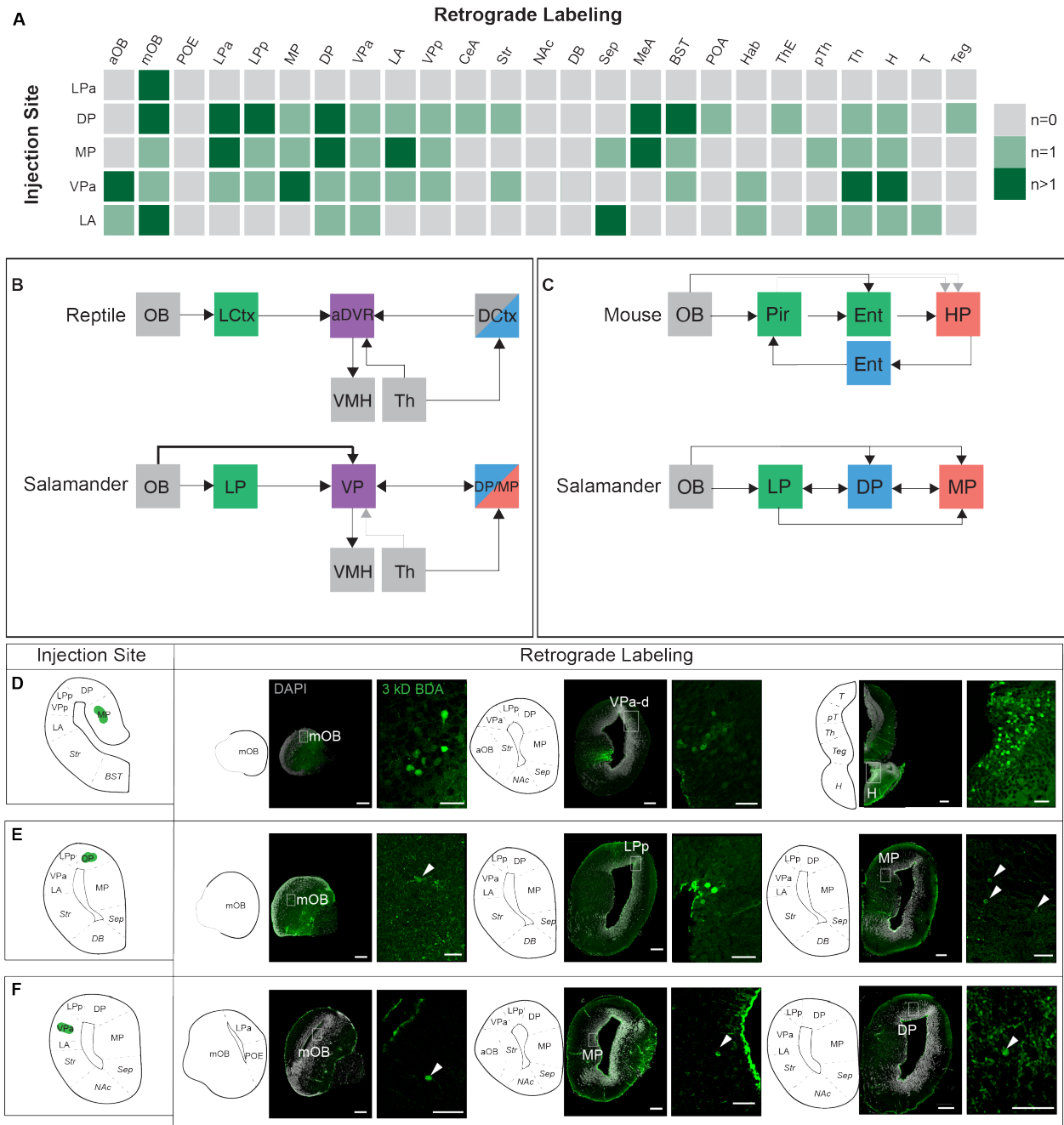


Fig. S19. Additional tracing data: injections of retrograde tracers in pallial regions. (A) Summary of regions with retrogradely labeled cells (columns) following 3 kD BDA tracer application into injection sites (rows): LPa (n=2), DP (n=2), MP (n=2), VPa (n= 4), LA (n=2). Connections are labeled in light or dark green when observed in one or more than one injected brain, respectively. **(B-C)** Schematic connectivity maps of salamander pallial regions, as compared to mouse (93) and reptile. Gray lines represent sparsely documented connectivity, thick lines represent well documented connectivity. **(D-F)** Additional representative images of 3 kD BDA tracer injections and retrograde labeling into **(D)** medial pallium; **(E)** dorsal pallium; **(F)** anterior ventral pallium. For neuroanatomical abbreviations, see Fig. S1. Scale bars: 200 μ m for slice overviews; 50 μ m for magnified images.

Other Supplementary Materials

Movie S1. Overview of the morphology of an entire *Pleurodeles waltl* brain. Optical coronal slices (1741 sections) of a cleared brain, stained with the nuclear marker TO-PRO-3, and imaged using light sheet microscopy along its anterior-posterior axis.

Movie S2. Expression of marker genes in entire brains in virtual sections. Optical coronal slices (~1800 sections per brain) of intact, cleared brains, stained using HCR or immunostaining, and imaged with light-sheet microscopy along their anterior-posterior axis. Slices of entire brains expressing *Etv1* (highly expressed in MP and DP), *SATB1* (highly expressed in LP), *Sox6* (highly expressed in POE and VPa/p), *Slc17a6* (highly expressed in mOB, VPP, and Amy), *Nr2f2* (highly expressed in VPP/MeA/LA), *Penk* (highly expressed in mOB, Str, and MP, among others), and *Rorb* (highly expressed in LP, VPa, and VP) are shown sequentially.

Movie S3. Expression of marker genes in entire brains in 3D. 3D maximum intensity projections in intact, cleared brains, stained using HCR or immunostaining, and imaged with light-sheet microscopy, each rotated along their medial axis. 3D brains expressing *Etv1* (highly expressed in MP and DP), *SATB1* (highly expressed in LP), *Sox6* (highly expressed in POE and VPa/p), *Slc17a6* (highly expressed in mOB, VPP, and Amy), *Nr2f2* (highly expressed in VPP/MeA/LA) and *Penk* (highly expressed in mOB, Str, and MP, among others) are shown sequentially.

Data S1. DNA sequences for colorimetric ISH and HCR ISH probes.