



Supplementary Materials for

The dynamics of gene expression in vertebrate embryogenesis at single-cell resolution

James A. Briggs, Caleb Weinreb, Daniel E. Wagner, Sean Megason,
Leonid Peshkin, Marc W. Kirschner,* Allon M. Klein*

*Corresponding author. Email: marc@hms.harvard.edu (M.W.K.);
allon_klein@hms.harvard.edu (A.M.K.)

Published 26 April 2018 on *Science* First Release
DOI: 10.1126/science.aar5780

This PDF file includes:

Materials and Methods
Figs. S1 to S18
Table S1
References

Other Supplementary Materials for this manuscript include the following:
(available at www.sciencemag.org/cgi/content/full/science.aar5780/DC1)

Movies S1 to S3
Data S1 to S7

Materials and Methods

Collection of Xenopus embryos

Xenopus tropicalis embryos were collected as described in (46): mating pairs were injected with 200U of HCG per female and 50U HCG per male and allowed to mate naturally overnight at 22C. Embryos were dejellied in 2% wt/vol cysteine the next morning, ~6hrs after injection, and sorted into batches of staged clutches pre-mid-blastula-transition (pre-MBT) and allowed to continue developing at 25C in 1X MMR. As the embryos reached the desired developmental stage, 5-10 embryos at a time were sampled from the clutch for dissociation and processing by InDrops. A single healthy clutch was used for the first InDrops time series; a further five independent clutches were used for the timepoints comprising the second replicate experiment.

Development of a dissociation media for Xenopus embryos

To identify the best possible dissociation conditions for *Xenopus* embryos prior to InDrops RNA sequencing, we initially performed a broad screen of known cell dissociation enzymes and buffers, including calcium magnesium free media (CMFM) (46), Newport dissociation media (2), trypsin-EDTA, TrypLE (ThermoFisher 12604-013), Accutase (Innovative Cell Technologies AT104), papain and collagenase. Five devitellinized embryos were incubated in each buffer at 18°C in 12 well plates, with gentle swirling by hand every minute initially, and then every 10 minutes until dissociation was complete or 1 hour had passed. All enzymatic dissociation reagents were ineffective on early *Xenopus* embryos. CMFM and Newport buffers led to substantial cell dissociation, but CMFM completely failed to dissociate the pigmented outer layer of cells. Newport buffer was deemed the best condition among those screened, but was still slow, requiring up to 1hr for near complete dissociation, and it frequently failed to completely dissociate small clusters of cells in each embryo.

Newport buffer had originally been formulated for the dissociation of pre-MBT *Xenopus* blastomeres. We attempted to optimize it further for post-MBT embryos. We systematically varied salt composition, pH, and buffering acid used. Newport buffer [0.1M sodium isethionate, 20mM sodium pyrophosphate, 20mM glucose, pH 9] could be made substantially more potent by raising the pH to 10.5 in the presence of 3-(Cyclohexylamino)-1-propanesulfonic acid (CAPS; a Good's buffer with pKa=10.4) (Sigma C2632). In this new buffer – which we term “Newport 2.0” [0.1M sodium isethionate, 20mM sodium pyrophosphate, 10mM CAPS, pH 10.5] – embryos completely dissociated in <30min nutating at 80rpm and 18C. Newport 2.0 effectively dissociated embryos from stage 8 (earlier stages not tried), through stage 22 (early organogenesis). After stage 22 certain tissues fail to dissociate completely, likely because of newly generated extracellular matrix. Dissociated cells were >95% viable for up to 2hrs at RT after washing with PBS-/- (Corning; 21-040-CV) as determined by staining with DAPI (wash protocol detailed below; DAPI added at 0.1ug/mL). We also tested the extent of cell lysis by measuring free floating RNA in solution after dissociation by fluorometry with RiboGreen dye (ThermoFisher). After dissociation with Newport 2.0, free RNA was <5% of total embryo RNA, as judged by measurements of RNA in solution following dissociation compared to dilution series of RNA from whole embryos. We did not systematically investigate why high pH improves dissociation; it is possible that the high

pH increases the dissociation constant of calcium ions bound to integrin receptors that hold embryonic cells together, or that it causes non-specific surface protein denaturation.

Dissociation of *Xenopus* embryos into single cells for InDrops scRNA-seq

Removal of vitelline membranes: Vitelline membranes were removed by treatment with 1mg/mL pronase (Calbiochem) in 1X MMR (2), nutating at 80rpm and 18C for 8-10min. Forceps and inspection on a dissecting microscope were used to judge the earliest point at which the vitelline membrane had dissolved to minimize damage to the superficial layer cells of the embryos. Embryos were washed 2-3 times gently in 1X MMR immediately following pronase treatment to remove any vitelline membrane fragments stuck to the surface of the embryo, since carry over vitelline membrane interfered with dissociation.

Dissociation: A 12 well plate was pre-treated with Bovine Serum Albumin (BSA), and then a single well was filled with 5 mL of Newport 2.0 dissociation buffer [0.1M sodium isethionate, 20mM sodium pyrophosphate, 10mM CAPS, pH to 10.5 using NaOH]. Washed, devitellinized embryos were then transferred into the well. BSA pre-treatment was performed by air drying a film of 1mg/mL UltraPure BSA (Thermo Fisher) in ddH₂O under a heat lamp. This treatment prevents dissociated cells from sticking to the tissue culture plastic surface. The 12 well plate well was filled to the top with Newport 2.0 buffer, and sealed using a layer of parafilm and silicone well isolator (Sigma; S0810) after excluding all air bubbles. Sealing from air was essential to prevent cell lysis at the air-liquid interface during mechanical agitation. To mechanically accelerate dissociation, the sealed well plate containing fresh embryos was taped shut and attached to a benchtop Vortex genie using a plate adaptor. Embryos were agitated on speed setting 1-3 for 5 min, then accelerated to speed setting 7 for an additional 5-10 min or until dissociation was complete (typically <20 min). Embryos were entirely dissociated into monodisperse cells following this procedure. As soon as dissociation was complete, the plate was cooled on ice and single cells were allowed to settle to the bottom of the well by gravity for 5-10 min.

Washes and cell collection: Once settled, the single cell suspensions were washed by letting the cells settle by gravity through fresh buffer underlying the cells. Cells were first washed through a cold solution of 5% v/v Optiprep (Sigma; D1556) in 1X PBS^{-/-} (Corning; 21-040-CV), then increasing to 10% v/v and 20% v/v. The buffers are added slowly against the bottom of the well, taking care to avoid cells and without agitating the cells so that they would not rise to the air-liquid interface. Each increasingly dense layer of Optiprep replaced the bottom liquid around the cells, allowing complete buffer exchange without the need for centrifugation. Cells require 1-2 minutes to settle after each wash. After the washes, 1mL of 30% v/v optiprep with 1mg/mL UltraPure BSA (Thermo Fisher) in PBS^{-/-} was added as a loading solution for the cells. 30-35% v/v Optiprep is approximately neutral density for *Xenopus* cells, preventing them from settling during an InDrops run. To load the cells into a syringe for the inDrop platform (1), cells were floated off the bottom surface of the well by adding 0.5mL 40% Optiprep (denser than the cells) in PBS^{-/-}, causing them to accumulate at the interface between the 30% and 40% Optiprep phases. The cell layer was then aspirated directly into a 1mL syringe, which was backfilled with 0.1mL loading solution to avoid contact between cells and air, minimize cell handling, and block the plastic syringe surface. A wide bore needle

tip and tubing (0.86mm internal diameter; SCI; BB31695-PE/5), pre-coated with BSA, was then attached to the syringe and connected to the InDrops microfluidic setup for cell injection. InDrops droplet collection was initiated immediately. The viability of dissociated cells remains >95% over 2hrs in loading solution.

InDrops single cell barcoding, library preparation, and sequencing

Dissociated *Xenopus* cells were processed by InDrops as described by Klein et al (*1*) with the following changes: we minimized the distance between the cell injection syringe and the microfluidic chip; used a wider bore size tubing; and coated both the tubing and syringe with BSA prior to performing the experiment. These modifications minimize shear stress during microfluidic processing which can badly damage the large and fragile *Xenopus* blastomeres. After reverse transcription in droplets, barcoded emulsions were split into batches estimated to contain 2,500 cells, and chemically broken. Libraries were prepared as described previously (*1*). For the initial replicate of 42,385 cells, single cell libraries were pooled and sequenced together across seven NextSeq high runs, generating a total of ~3billion raw paired-end reads (**Table S1**). The second replicate of a further 94,131 cells was sequenced across a total of ten NextSeq high runs achieving approximately half the sequencing depth per cell of replicate 1 (**Table S1**). Within each replicate pool concentrations were adjusted to ensure similar depth across all samples. Reads were demultiplexed using an updated version of the custom bioinformatics pipeline described in Klein et al (*1*) available at <https://github.com/indrops/indrops>, using a reference *Xenopus* transcriptome (XTr 9.0). The final output is an integer counts matrix of cells vs. genes.

Expansion of gene symbol assignments in the XTr 9.0 reference transcriptome

Over 15k genes in the XTr 9.0 reference transcriptome lack gene symbols, though in many cases an unambiguous protein identity can be identified by sequence homology. To fill holes in the gene symbol assignments we assigned protein gene symbols to each XTr 9.0 transcript using a modified reciprocal best HMMER hit approach (*47*) based on a target reference set of curated human proteins. Specifically, as a reference set we used 25,208 Human and 4,492 *Xenopus tropicalis*, and 6,430 *Xenopus laevis* curated sequences obtained from UniProt (*48*). Multiple matches with identical E value were resolved based on bit score. Additionally, sequences which did not fall into any bidirectional best match were assigned a symbol based on the unidirectional best HMMER match. Only matches with an E value better than $1e-5$ were considered. The result of our pipeline matched a total of 22,599 *Xenopus tropicalis* genes to one of 16,016 unique gene symbols, which provides 11,341 more *Xenopus* genes with a gene symbol than available in the original genome annotation, significantly enhancing data usability. We provide the enriched gene symbol assignments alongside the XTr 9.0 gene symbol assignments on our web-browser and in **Additional Data Table S7**.

Data clean up: minimum expression threshold, doublet and background removal

To remove background signal from putative empty droplets, only cells with >500-1000 detected UMIs, depending on the library sequencing depth (see **Table S1**), were carried forward for analysis. After plotting cells in two dimensions using tSNE (details below), we observed some clusters that appeared to consist of cell doublets and others

that likely represented empty droplets. We classified doublet clusters using three criteria: i) lack of specific marker genes; ii) mixing of marker genes from other clusters; iii) a small size matching the experimental expectation of <2.5% double cell encapsulation during inDrops. Cells associated with doublet clusters were excluded. Background clusters were distinguished by i) lack of specific marker genes; ii) weak and uniform expression of all genes observed in the time-point. We noted that these background clusters form a characteristic ball-like shape in tSNE. All analyses presented in the paper represent cells passing these three filters (UMI>minCounts, not-doublet, not-background). We provide unfiltered gene expression data online and indicate which cells were excluded.

Visualization of data using PCA-tSNE and SPRING

To visualize high-dimensional single cell data in **Fig. 1** we implemented the PCA-tSNE pipeline described in Klein et al (*1*), which follows these steps: i) normalize gene expression by the total transcript count per cell; ii) identify principally variable (PV) genes and non-trivial principal components compared to randomized data; iii) z-score normalize PV genes, perform principal component analysis (PCA) and retain the non-trivial principal components; (iv) run tSNE (*49*) to generate a two-dimensional visualization. In addition to tSNE, we used SPRING as a complementary method to visualize single cell data (e.g. in **Fig. S6**), as described in (*50*). The tSNE perplexity parameter was set to 25 in all plots. SPRING builds a k-nearest neighbor (knn) graph from the high dimensional distances between cells and then embeds the graph in two dimensions using a force-directed layout, which emphasizes the continuous relationships between different cell clusters. For SPRING we used Euclidean distances computed from the same PCA coordinates used as inputs to tSNE.

Clustering of cell states

We first clustered cells from each time-point using local-density clustering (DBSCAN, as used in (*33*)) applied to tSNE coordinates. We then used two steps to refine our clustering. First, we performed an initial round of clustering on all cells that over fragmented the data and then recombined adjacent clusters that had less than 3 genes expressed at >0.5UMIs/cell on average that were >4-fold enriched in one state versus the other. Second, we isolated every cluster and performed a second round of dimensionality reduction and clustering on the cells for each cluster in isolation. The clustering approach can be summarized as follows: i) apply PCA-tSNE to the desired normalized counts matrix (see above); ii) use DBSCAN to generate an over-fragmented clustering; iii) manually merge adjacent clusters with insufficient differential gene expression – we refer to the clusters generated in steps (i-iii) as “level 1” clusters; iv) repeat steps (i-iii) for each level 1 cluster to generate sub-clusters, which we refer to as “level 2”. This process was performed on replicate 1 cells from stages 8 – 22 (all even stages), and on new timepoints from replicate 2 – stage 11 and 13. Other replicate 2 cells were clustered by a method that was informed by the initial clustering (see section below “Addition of replicate cells to cell state tree scaffold”). Three rare cell states were classified by manual inspection of marker genes as they were too few in number to form a distinct cluster recognizable by tSNE-DBSCAN: germ cells, neuroendocrine cells past S14, and hatching gland cells past S16.

Benchmarking against an alternative clustering method

We tested the robustness of tSNE-DBSCAN clustering by comparing clusters of stage 22 (replicate 1) cells to those generated by an alternative method based on spectral clustering with a kNN-graph kernel. We used the graph generated by SPRING, described above. Notably, the spectral clustering approach depends directly on the high dimensional distances between cells as opposed to their 2D tSNE coordinates. We found good agreement between the spectral clusters and those generated through our semi-manual tSNE-DBSCAN approach (**Fig. S5**). Using an equal number of spectral and tSNE-DBSCAN clusters, a median of 79% cells in each tSNE-DBSCAN cluster mapped to its most similar spectral cluster, while a median of 89% cells in each spectral cluster mapped to its most similar tSNE-DBSCAN cluster. Most of the disagreements involved a single tSNE-DBSCAN cluster mapping to multiple spectral clusters (or vice versa), rather than a violation of cluster boundaries.

Annotation of cell state clusters

To relate the scRNA-seq datasets to known embryonic cell types we initially manually matched each gene expression cluster to Xenopus Anatomy Ontology (XAO) (8, 9) terms based on top marker genes. The XAO includes thousands of curated *in situ* staining measurements of tissue specific gene expression. By identifying genes that were specific to each gene expression state, and that had clear tissue specific expression patterns on XAO, it was possible to associate each gene expression state to a particular cell type (also see main text). We then validated our annotations using a systematic approach as follows. (1) For each XAO term, we bioinformatically curated a set of literature-validated marker genes by parsing the gene expression database for *Xenopus laevis* and *Xenopus tropicalis* on Xenbase, found at the following locations:

ftp.xenbase.org/pub/GenePageReports/GeneExpression_tropicalis.txt
ftp.xenbase.org/pub/GenePageReports/GeneExpression_laevis.txt

We filtered for genes measured by RNA in situ hybridization only. (2) For each annotated single cell cluster, we identified marker genes as described below in the section “Identification of marker genes for each cell state”. (3) For each XAO term, we intersected the marker gene list from (1) with the cluster-specific gene list from (2). From these automatic marker gene annotations, all but nine states had at least one Xenbase-shared marker gene, of which five (dorsal lateral plate region, intermediate mesoderm – *ssg1*, trigeminal and profundal placodes, small secretory cells, and epidermal progenitor – *tp63/tl12*) had matched anatomy terms with marker genes documented on Xenbase in the first place. 54% of annotations had at least three marker genes (compared to 5% when tissue annotations were randomized). The resulting list of literature-supported marker genes for each cluster is provided in **Additional Data Table S1**, along with literature references.

Connecting cell states across time from gene expression similarity

We investigated the developmental transitions connecting clusters across time points by asking each cluster to ‘vote’ on its most likely ancestor cluster from the previous time-point, using the following steps.

1. For each two adjacent time points, we embedded all cells from the two time-points in the PCA space learned from the second time point only, keeping non-trivial PCs as defined above. This embedding causes cell-cell distances to reflect gene expression variation between tissues, as opposed to global changes over time.
2. In this embedding, for each cluster in the late time point, each constituent cell reported on the cluster identity of its 5 nearest neighbors from the previous time point using a Euclidean distance metric.
3. The number of edges to each early time point cluster were aggregated across all cells in each cluster. Percentages of votes cast for each possible ancestor are shown in heat maps on **Figure S7**, and the cell state tree presented in the main text represents vote winners. The winner of each vote was usually unambiguous, with <5% of assignments sharing more than 30% of votes with non-ancestor states, and a median vote-winning share of 88%.

We performed the ancestor voting process first for the “level 1” clusters (see section on clustering for a definition) and then, for level 1 clusters that had been sub-divided into level 2 clusters, we performed a second – private – voting process on the level 2 clusters – i.e. between only subpopulations contained within connected level 1 states. This two-step procedure was necessary because the level 2 clusters are only distinguishable with PCA coordinates learned from each level 1 cluster separately, and therefore had to be bioinformatically isolated when computing inter-time point cell-cell distances.

Addition of replicate cells to cell state tree scaffold

Late during the preparation of this manuscript, we performed a large replicate experiment, that added a total of 94,131 additional single cell transcriptomes to the original time series. The replicate included all original timepoints except stage 10. It also included two new intermediate timepoints – S11 and S13. All cells were quality filtered by the same process as the original experiment. New timepoints (S11 and S13 only) were clustered and mapped into the cell state tree structure by the same methods as the original experiment.

New data from existing timepoints (all but S11 and S13) were clustered by a distinct process that (i) first assigned every cell to a pre-existing annotated cluster; (ii) performed a batch correction of gene expression between the biological replicate experiments; and then (iii) examined the combined data set for novel sub-structure within the annotated clusters. The final annotations presented in this paper incorporate new sub-structure found through this process. The detailed steps are as follows:

- (i) *Assign cells to a pre-existing annotated cluster:*
 - a. New data at each time point was fragmented into $4*N$ clusters, where N is the number of clusters in original experiment at the matched time point. Clustering was performed by k-means spectral clustering on the kNN graph, constructed using parameters described above (see SPRING visualization).
 - b. Each cell in the new data was assigned a cluster label based on the original cluster annotations, by connecting the $4*N$ new clusters to the original N clusters within each time point. The assignment algorithm is as described for cross-timepoint mapping (See previous Methods section; median voting consensus of 80%).

- (ii) *Batch correct biological replicates*: This method is a variation on the batch correction algorithm recently proposed by Marioni et al. in pre-print (<https://doi.org/10.1101/165118>). From step (i), each cell in both biological replicates is assigned to a unique cluster i . We defined the gene expression vector for cell j in cluster i from replicate 1 as x_{ij} , and respectively from replicate 2 as y_{ij} . The centroids of cluster i in the two replicates (average of all cells in the cluster) x_i, y_i respectively, were calculated in units of TPMs from raw counts. We further defined the arithmetic mean of both centroids $z_i = (x_i + y_i)/2$, and a multiplicative correction factor $\Delta_i = \frac{1+x_i}{1+z_i} - 1$. The corrected gene expression values for each cell j in cluster i were then calculated as $x'_{ij} = \frac{1+x_{ij}}{1+\Delta_i} - 1$, and $y'_{ij} = \frac{1+y_{ij}}{1-\Delta_i} - 1$. Corrected counts matrices x' and y' from both experiments were combined and carried forward for further analysis. This correction procedure centers the centroids of cells in each cluster from the two biological replicates through a multiplicative correction, ensuring non-negative and non-infinite values for gene expression. After correction, both replicates intermixed completely.
- (iii) *Discover substructure within pre-annotated clusters*: each new cluster, now composed of cells from both replicates, was sub-clustered in isolation by tSNE-DBSCAN (see above). If only a single cluster was identified, then no sub-structure was added to the annotation. If more than one cluster was identified, the new sub-clusters were examined for expression of marker genes of the original annotated cluster. If a sub-cluster lacks typical marker genes of the tissue, the cluster was excluded as outlier cells. All sub-clusters expressing marker genes of the original annotated cluster were assigned novel annotations and included in the final state tree.

Benchmarking of cell state tree similarity relationships against XAO

We manually inspected each of the 257 edges in the resulting cell state tree, using the XAO database and literature to determine the validity of each edge. Each edge in the tree links a parent cluster in time point t to a daughter cluster in time point $t+1$. We tested three criteria: (a) the daughter cell state is believed to arise from the parent cell state; and (b) if the daughter cell state differs in XAO annotation from the parent, then the daughter state should appear at the time specified in the XAO database indicated in the Start_Stage field; (c) if the daughter cell state differs in XAO annotation from the parent, then the parent state should end at the time specified in the XAO database indicated in the End_Stage field. For each edge, we provided at least one reference to XAO or to the literature supporting our verdict. We assigned each edge one of seven flags: “OK” indicating that all criteria are satisfied (203 edges); “Discovery: early” indicating that a state appears prior to its annotated first appearance in the XAO (18 edges); “Starts too late” indicating that a state only appears after its annotated first appearance in the XAO (9 edges); “Ends too early” indicating that the parent state terminates prematurely (1 edges); “Ends too late” indicating that the parent state terminates after expected in the XAO (3 edges); and “Error” indicating that the daughter does not arise from the parent state (2 edges). In addition, 21 edges were assigned “N/A” as they related to the appearance of daughter states that did not correspond to XAO terms. The full table is

provided in **Additional Data Table S2** and the difference in timing of appearance is plotted in **Fig. 3A**.

Identification of marker genes for each cell state

The 259 clusters in the data were mapped to 87 unique annotations in the *Xenopus* cell state tree, many appearing over multiple time points (also see above section “Annotation of cell state clusters”, and main text). We defined marker genes for each of the 87 unique annotations in the *Xenopus* cell state tree with the following criteria: i) average expression enriched >5-fold compared to the average for the rest of embryo, ii) expression detected in >15% of cells in the state, and iii) the state has the highest average expression for the gene. We identified 2,159 marker genes across all states, with a median of 6 marker genes per state. The mean enrichment for the most specific marker gene for each state was 93-fold, and the mean expression level across all identified marker genes was 2.2 UMIs (930 TPM), a value approximately half the average expression level of the ‘housekeeping gene’ beta-actin (5.65 UMIs/cell or 2,360 TPM). Marker gene lists are tabulated in **Additional Data Table S4**.

Differential gene expression at cell fate choices

Differentially expressed genes between two cell states were identified using two-tailed Wilcoxon Rank Sum Tests, at a stringency of $p < 0.001$ with multiple hypothesis testing correction and FDR of 5%. Before rank sum tests were performed data was total count normalized across states, and filtered for genes that were expressed in ≥ 10 cells, had an average of ≥ 1 UMI/cell in at least one state, and were ≥ 4 -fold differentially expressed between states. Overall this procedure is highly conservative, requiring both high expression and high fold-change in expression; it identifies the most prominent differences between states efficiently and with a low false positive rate. Differentially expressed genes at each fate choice are tabulated in **Additional Data Table S3**.

Cluster continuity analysis

We used the kurtosis of the cell distance distribution as a metric of cell clustering (**Fig. S4**). Specifically, the metric c_i is calculated for each cell i , at each timepoint, as the $\ln(\text{kurtosis})$ of the distribution of Euclidean distances from cell i to every other cell at the same timepoint, in the principal component latent space defined above in section “Visualization of data using PCA-tSNE and SPRING”. The metric serves as a continuity index because the appearance of discrete clusters leads to a bimodal or multimodal distance distribution for each cell, which increases the kurtosis. This analysis was performed on replicate 1 data.

Alignment of *Xenopus* and Zebrafish cell state trees

Xenopus and Zebrafish cell states were matched manually to enable computational comparisons of gene expression across species. Many-to-many matches of cell states were allowed to avoid ambiguities in matching timepoints across species, and to allow groups of states showing greater substructure in one species that the other to nonetheless be compared (e.g. the zebrafish spinal cord was divided into several subpopulations compared to the single *Xenopus* state). The criteria to match states were that they should: i) share a common position within the developmental hierarchy of both species according

to Zfin and Xenbase anatomy ontologies, ii) share at least 3 validated marker genes, and iii) appear at approximately the same time across species. We sometimes encountered matchings where: 1) the name used by the frog community (Xenbase) differs slightly from that used by the zebrafish community (zfin); or 2) the time point at which we profiled the two species lead to one corresponding tissue having a less mature name as compared to the other (despite highly similar gene expression). To bioinformatically validate the manual matching we randomly permuted the cell state matchings and asked whether any changes could improve the mean gene expression correlation between species by >0.1 , indicating a potential error. This revealed no errors in the manual matchings.

Transcription factor (TF) re-use analysis

To identify TFs that were activated multiple times during early development, we used a graph based approach to search for TFs with multiple non-overlapping expression domains on the cell state tree. The cell state tree represents a directed graph. The general approach is (1) to binarize the gene expression for each gene on the cell state tree; and (2), to count the number of connected components of the graph after discarding all nodes that are non-expressing. This simple approach is elaborated with additional steps in order to (a) merge connected components that are close on the cell state tree, to conservatively allow for noise-induced “Gene-off” events; (b) filter out connected components consisting of just a single state or collections of isolated states due to merging, as these may be due to noise. We analyzed the list of differentially expressed (DE) TFs (the intersection of a TF list from (51) and the DE gene list described above, filtered to exclude broadly expressed genes defined by a skewness over cluster averages less than 2). Then, for each gene, the detailed steps are as follows: (1) Gene expression is binarized on the cell state tree, classifying a cell state as ‘gene-on’ for a given gene if: i) the average expression was $>x$ UMI/cell; ii) $>y\%$ of cells had detectable expression. (2) Additional edges are added to the cell state tree, connecting each cell state to all other cell states within z fate splits apart, i.e. states that can be accessed by walking on the graph and crossing no more than $z-1$ nodes that have >2 outgoing edges. (3) All “gene-off” states and their associated edges are discarded from the tree, to generate a graph consisting of “gene-on” cell states only. (4) A connected component of the “gene-on” graph was defined as robust if more than 50% of nodes in the component had at least one neighbor on the original cell state tree that is ‘gene-on’. Connected components that were not robust were discarded. (5) The number of connected components is the number of times a gene is activated on the cell state tree. The parameters x,y,z determine the sensitivity at which a gene is considered expressed, and the resolution of domains of expression on the cell state graph. We chose values for these parameters that maximized the accuracy of calling a TF as single-use or re-used for a “ground truth” training set of 20 TFs that were manually inspected to be expressed only once, or more than once. The values that achieved maximum accuracy (18/20 TFs called correctly) were $x=0.1$ UMI/cell (~ 500 transcripts per million (TPM)); $y=15\%$ cells; $z=3$ fate splits. The re-used TF analysis was performed on replicate 1 cells embedded in version 1 cell state tree. Results were tabulated in **Additional Data Table S5**.

Multi-lineage priming analysis

To systematically identify occurrences of multi-lineage priming (MLP), we searched for genes that were robustly co-expressed in the same cells at one point in time, and then still expressed but in disjoint sets of cells at a later point in time. Further, the latter cells are required to be descendants of the former cells in the cell state tree. To this end, we first classified pairs of differentially expressed genes in each of the fate splits on the tree as MLP+ or MLP-. Each individual gene was then deemed MLP+ with respect to a given fate split, if it contributed to at least one MLP+ pair at the fate split. Each gene pair was classified as MLP+ at a fate split as follows. For each sub-tree, we define a co-expression index for each time point, as the number of cells co-expressing both genes at the time point, divided by the number of cells expressing either gene. A gene pair was then classified as MLP+ for the fate split if: i) the maximum value of the co-expression index was x times higher than the minimum value for the subtree, with the max happening earlier in time than the min; ii) the maximum value was higher than y ; iii) the minimum value was lower than z . For *Xenopus*, we used $x=4$, $y=0.1$, $z=0.05$. For Zebrafish, we used the same parameters except $y=0.2$, to account for higher depth of sequencing. Note that the co-expression index is highly sensitive to the sequencing depth of the scRNA-Seq data, through the expression drop-out rate. Therefore, the absolute %MLP+ genes is an underestimate of the true value, but the qualitative trend in the %MLP+ genes over time is reliable, given a constant sequencing depth. The main result reported in the main text – a global trend of decreasing MLP over time – was insensitive to 2-3 fold changes in each parameter. Results were tabulated in **Additional Data Table S6**.

kNN-graph visualization of entire scRNA-seq timecourses

Whole-embryo kNN-graphs were generated using the approach reported in Wagner et al (sister paper). In **Fig. S9** we compare the *Xenopus* and Zebrafish cell state trees to this alternative representation.

Gene expression correlation between species

Gene expression conservation between species was measured using spearman correlation of the average normalized UMI counts for orthologous across 33 orthologous tissues. Each orthologous tissue contains one or more clusters from *Xenopus* and Zebrafish respectively, and the selection of clusters was performed manually with attention to marker gene expression, name, position in the embryo and inferred potential to produce downstream lineages. We created a list of orthologous genes and their percent sequence identities by selecting reciprocal best hits from a mutual amino acid BLAST search between all Zebrafish genes (assembly GRCz10) and all *Xenopus* genes (assembly XTr 9.0). Of these reciprocal best hits, only a subset consisted of genes that were robustly expressed in the single data. Reasoning that cross-species correlation of gene expression cannot be expected to exceed within-species correlation between different replicates, we computed a 'self-correlation' score for each gene by performing 50 random splits of the data from each species and recording the correlation of cluster-averaged expression for each gene in one half to itself in the other half. Only orthologous gene pairs where both genes had a self-correlation >0.75 were carried forward for analysis.

Analysis of functional gene categories enrich for conserved genes

Annotations of gene function, obtained from Gene Ontology (GO) term associations reported on Xenbase, were ranked by the proportion of constituent genes having a cross-species correlation of >0.66 , with this threshold chosen to establish an FDR of 1% against a null distribution generated from randomly matched (i.e. non-orthologous) genes. We restricted the analysis to orthologous gene pairs having a 'self-correlation' greater than 0.75, as explained above, and focused on GO terms with at least 20 genes among this set. P-values were assuming a binomial null distribution for the number of conserved genes among those associated with each GO term.

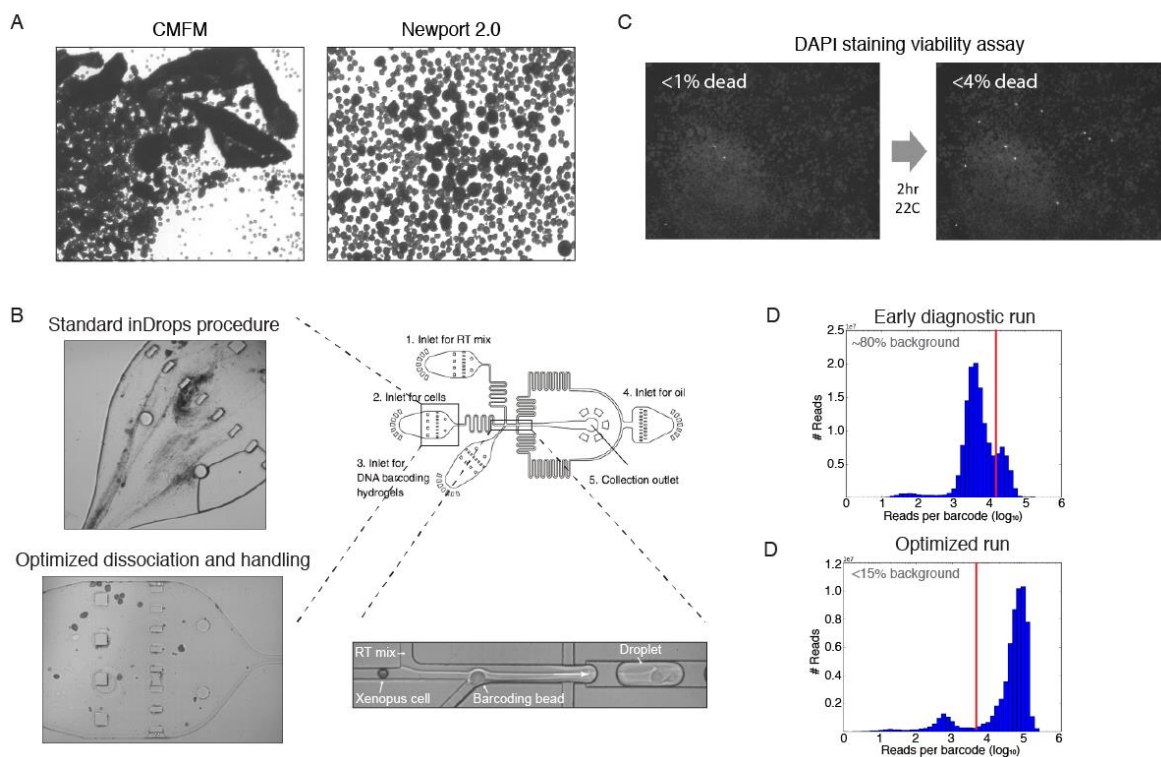


Fig. S1. Optimization of dissociation and inDrops procedures for Xenopus.

(A) Representative dissociation outcomes using either calcium magnesium free medium (CMFM) as compared to Newport 2.0 dissociation buffer. While CMFM fails to dissociate pigmented animal cap tissue, Newport 2.0 leads to complete dissociation. (B) Representative micrographs of *Xenopus* blastomeres during microfluidic droplet capture (also see Supplementary movies 1-3). Without optimizing cell handling procedures, lysis of cells during device injection led to severe chip blockages. (C) Cells dissociated using Newport 2.0 are >95% viable for over 2hrs at 22C. Arrows indicate DAPI positive dead cells immediately after dissociation and buffer exchange, and two hours later. (D) Reads per cell barcode histograms; right peak indicates signal coming from cells while left peak indicates background signal from ‘empty’ droplets. With optimized dissociation and handling procedures scRNA-seq runs with <15% total background signal, and <1.5% background signal per cell containing droplet were achieved.

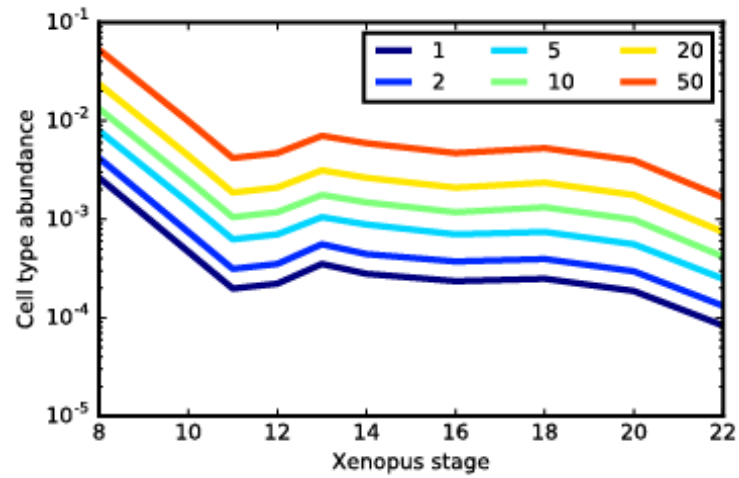


Fig. S2. scRNA-seq sensitivity analysis.

Cell type abundance (measured as proportion of cells in the embryo) that would allow the detection of N cells with 95% confidence, where N is indicated by line color, as determined by binomial sampling statistics. Values are given as a function of Xenopus stage, accounting for the different numbers of cells were profiled at each timepoint.

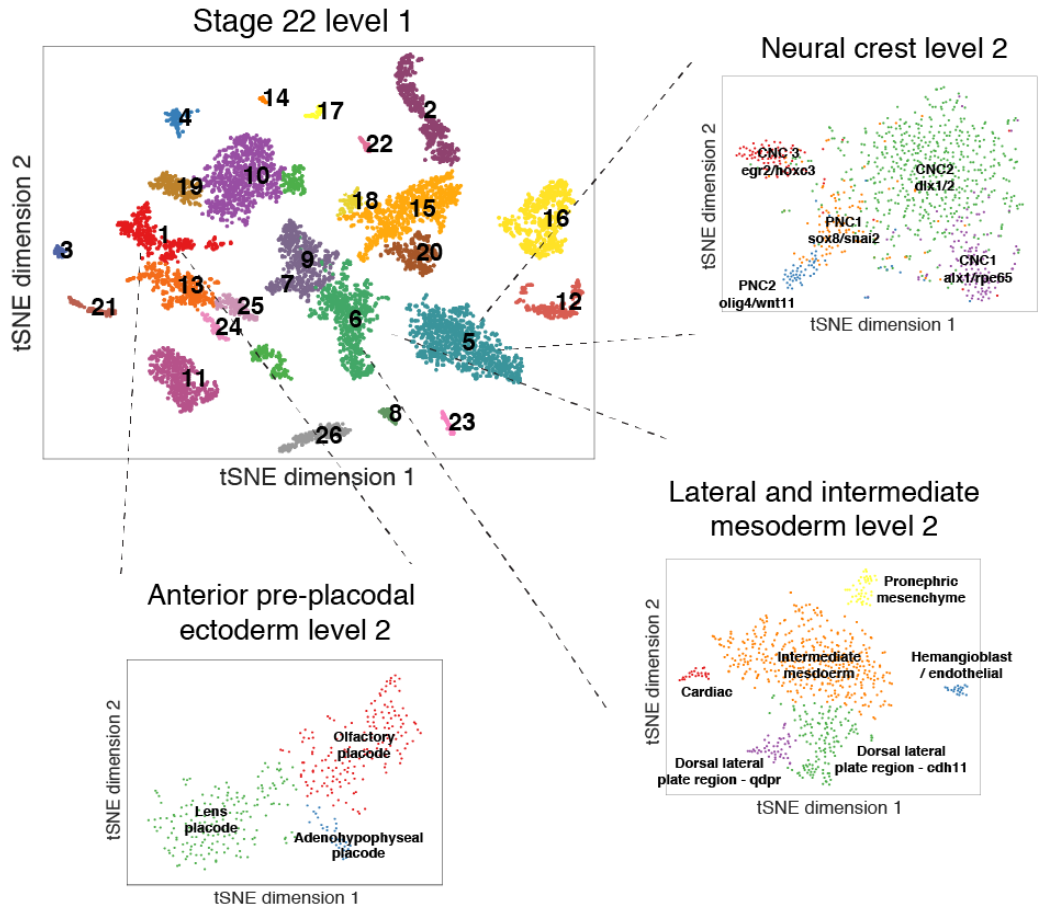


Fig. S3. Example of nested clustering.

To classify cell states, each timepoint was clustered in two rounds, or levels. The top left panel shows a representative level 1 clustering. The three breakout panels show how a nested round of tSNE-DBSCAN clustering on each individual cluster can reveal additional substructure, by focusing on variable gene expression within each subpopulation. Level 1 clustering key: 1 = anterior pre-placodal ectoderm; 2 = somite; 3 = ciliated epidermal progenitor; 4 = goblet cell; 5 = neural crest; 6 = lateral and intermediate mesoderm; 7 = neurons; 8 = lens progenitor; 9 = ventral mesoderm; 10 = neural plate; 11 = eye primordium; 12 = endoderm; 13 = posterior pre-placodal ectoderm; 14 = cement gland primordium; 15 = presomitic mesoderm; 16 = epidermal; 17 = ionocyte; 18 = tail bud; 19 = forebrain progenitor; 20 = intermediate mesoderm - *ssg1*; 21 = blood; 22 = notochord; 23 = migrating myeloid progenitors; 24 = cranial neuron; 25 = otic placode; 26 = small secretory cells. CNC = cranial neural crest; PNC = posterior (chordal) neural crest.

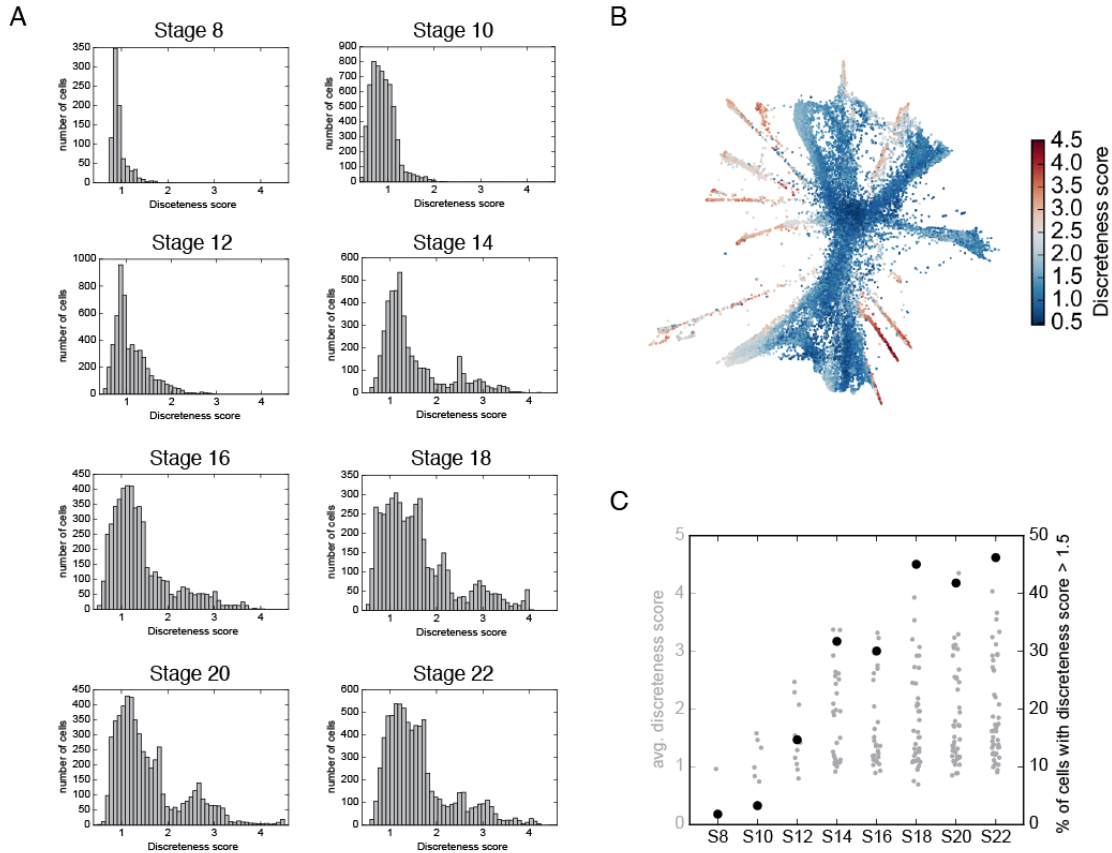


Fig. S4. Dynamics of discrete cell state appearance in *Xenopus* in scRNA-seq timecourse.

(A) For each cell, a ‘discretteness score’ was calculated as the log-kurtosis of its distance distribution to all other cells in gene expression space (see Materials and Methods section “**Cluster continuity analysis**” for further details). The premise of the score is that high kurtosis can indicate bimodality, or discrete clustering of the distance distributions. (B) Discretteness scores plotted on the all-cell kNN-graph from Fig. 2B, showing that scores are high on relatively discrete cell state projections. (C) Dynamics of discrete cluster emergence. Grey dots indicate the average discretteness score of each cell state at each timepoint. Discrete states with scores > 2 (higher than the right tail of the stage 8 distribution) first appear at stage 12 and accumulate over time. Continua are also apparent at each timepoint. The large black dots indicate the fraction of total cells at each timepoint with discretteness scores > 1.5. Scores shown are calculated on replicate 1 data.

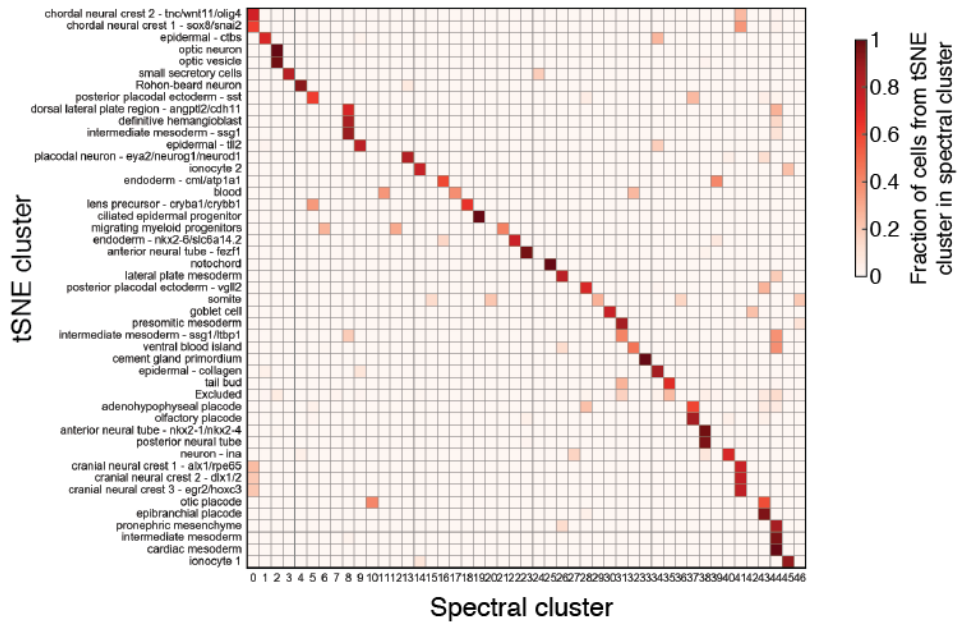
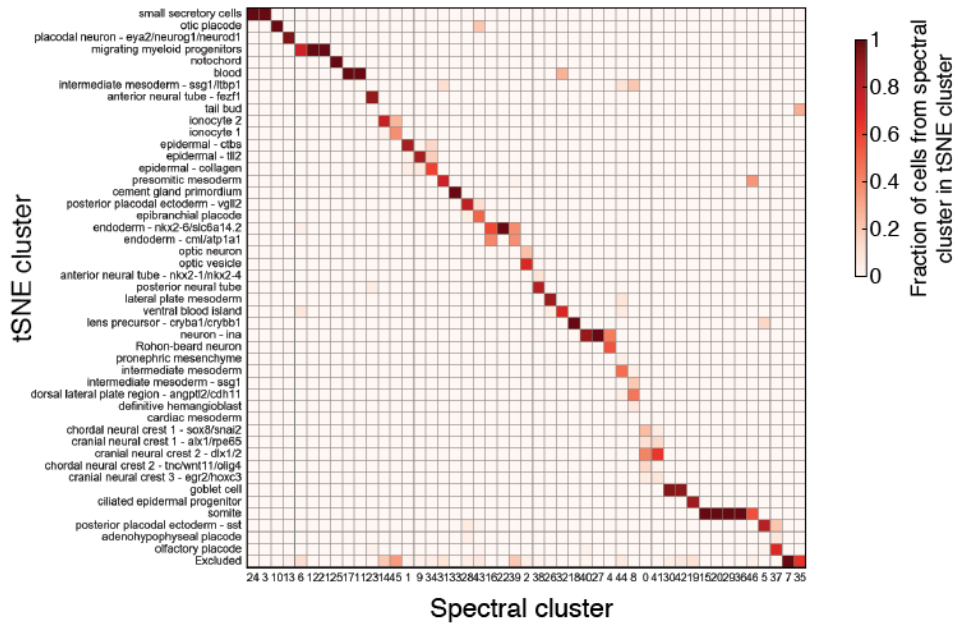


Fig. S5. Computational cross-validation of stage 22 tSNE-DBSCAN clustering by comparison to SPRING-spectral clustering.

To assess the robustness of the cell state classifications generated by tSNE-DBSCAN, we reclustered stage 22 scRNA-seq data (replicate 1) by an independent method – spectral clustering (see methods) – and compared the two results. Good qualitative agreement between the two approaches was observed: a median of 88.9% of cells in each spectral cluster mapped to its single most similar tSNE-DBSCAN cluster (top), while a median of 79% of cells in each DBSCAN cluster mapped to its most similar spectral cluster

(bottom). Occasionally a single tSNE-DBSCAN cluster mapped to multiple spectral clusters (or vice versa), but these disagreements in general did not indicate gross misclassifications. Off-diagonal signal generally involves <20% of cells in a cluster, and reflects subtle differences in cluster boundaries between methods.

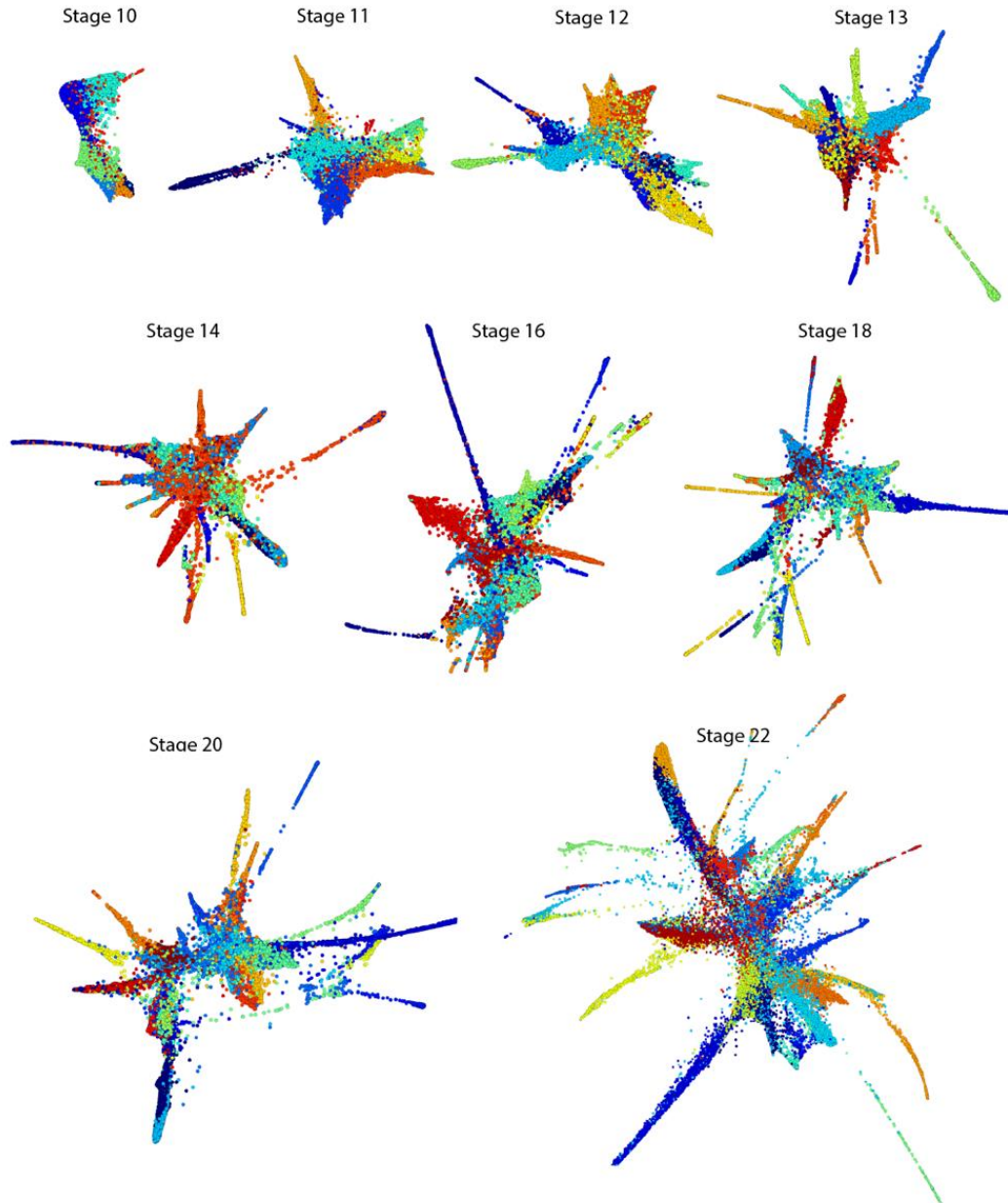


Fig. S6. Structure of tSNE-DBSCAN clusters is preserved on an independent visualization of the underlying scRNA-seq data.

To demonstrate that the gene expression structure underlying our cell state classifications across all timepoints was robust to different data processing methods (and therefore robust), we examined the structure of the clustering summarized in **Fig. 1B** on data visualizations generated by an independent method – SPRING. The plots show that the main text cell state classifications, defined using tSNE-DBSCAN clustering, and indicated by coloring here, retain their clustered structure on independent SPRING visualization of each timepoint. An interactive version of the SPRING plot for each timepoint with full annotations is available online at: tinyurl.com/scXen2018.

A — Main cluster mapping ————— B — Sub-cluster mapping —

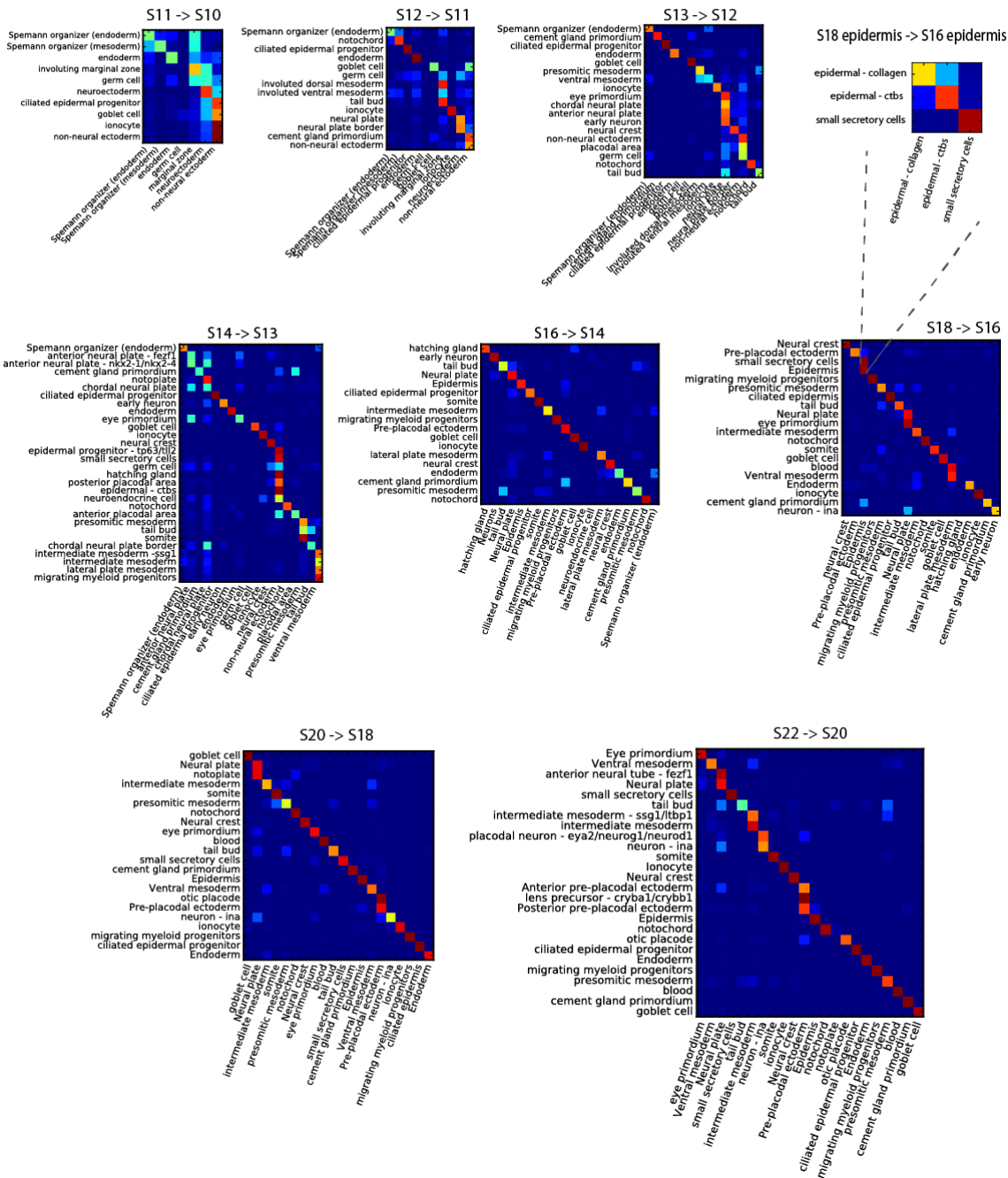


Fig. S7. Voting outcomes during ancestor assignment across timepoints.

(A) Voting outcomes during mappings between level 1 clusters. The median consensus was >88% across all assignments at all time-points. (B) Representative level 2 mapping of subpopulations within the S18 and S16 epidermis.

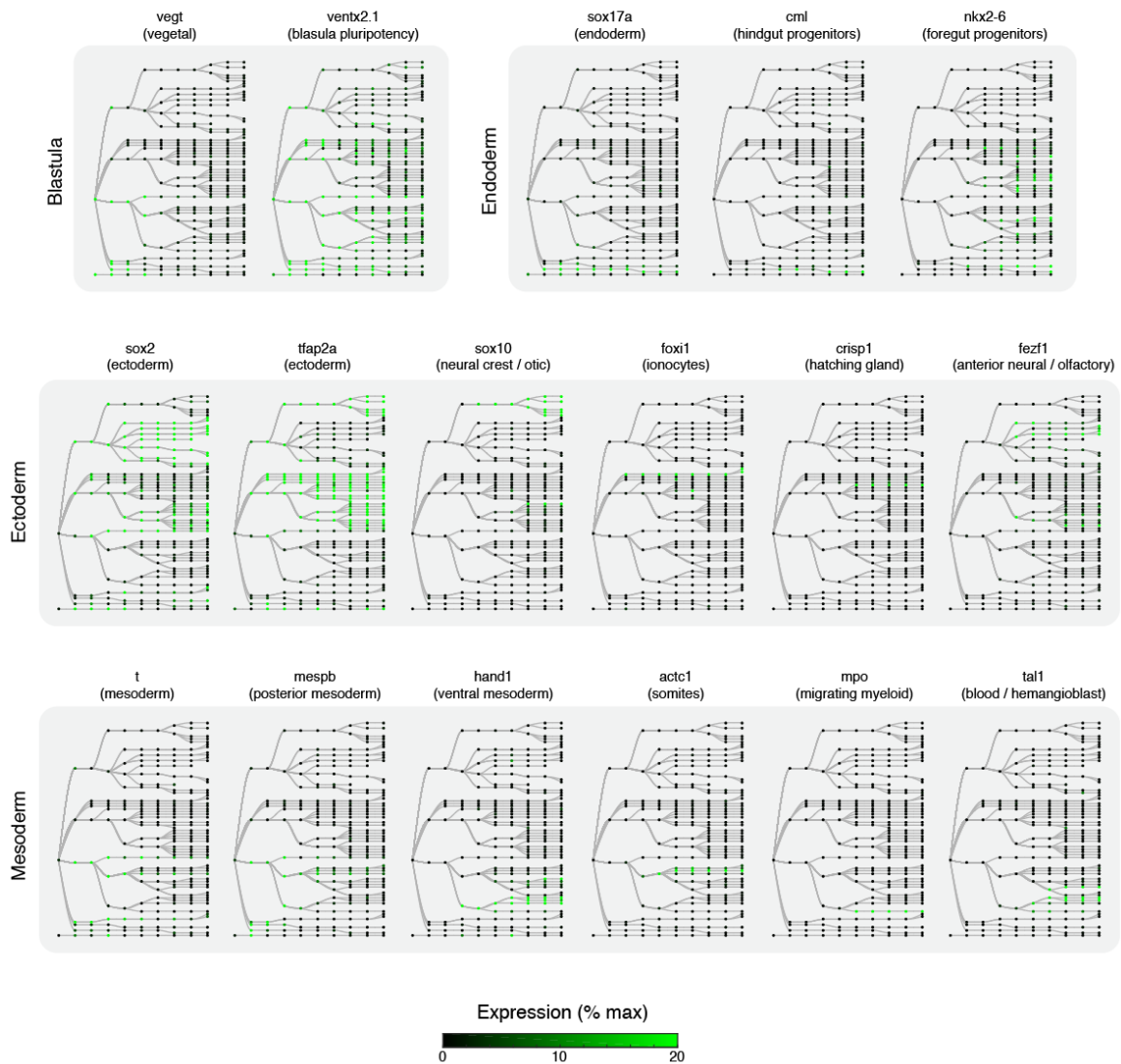


Fig. S8. Marker gene expression on the cell state tree.

Plots show the average expression level per node of the indicated marker genes on the cell state tree from main text Fig. 2C. Examples illustrate representative patterns for broad germ layer markers such as *sox2*, *t* (brachyury), and *sox17*.

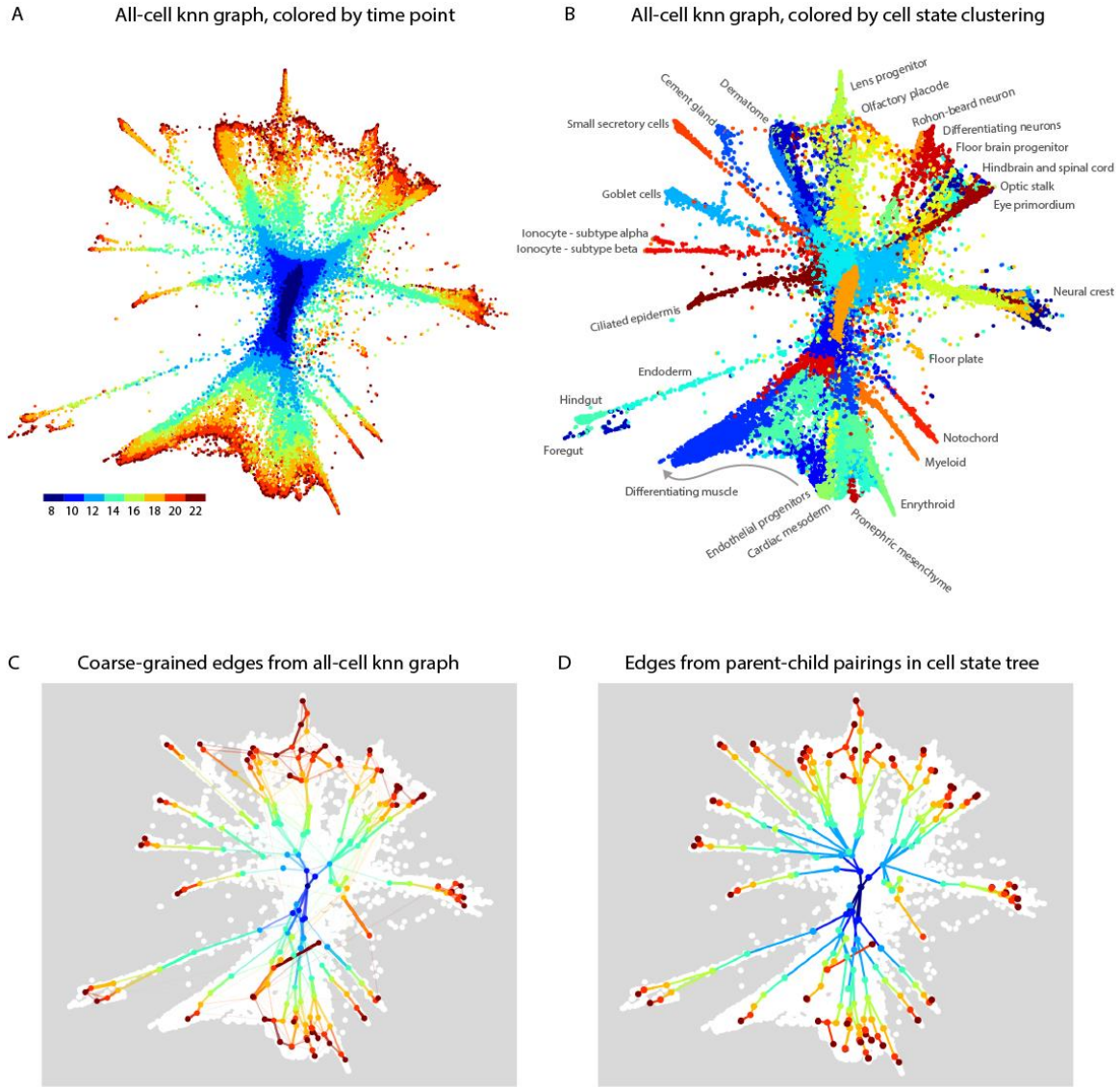


Fig. S9. Comparison of cell state tree structure to alternative cell state coarse graining algorithm.

(A) All cell kNN-graph from Fig. 2B showing the full *Xenopus* scRNA-seq timeseries (replicate 1; see **Table S1**) colored by timepoint. kNN graph was built using the algorithm presented in our sister paper (Wagner et al.). (B) All cell kNN-graph indicating cell state cluster identities. (C-D) Edge assignments between states across timepoints generated by coarse graining of the underlying kNN-graph (C) closely resemble those generated by the cell state tree algorithm (D).

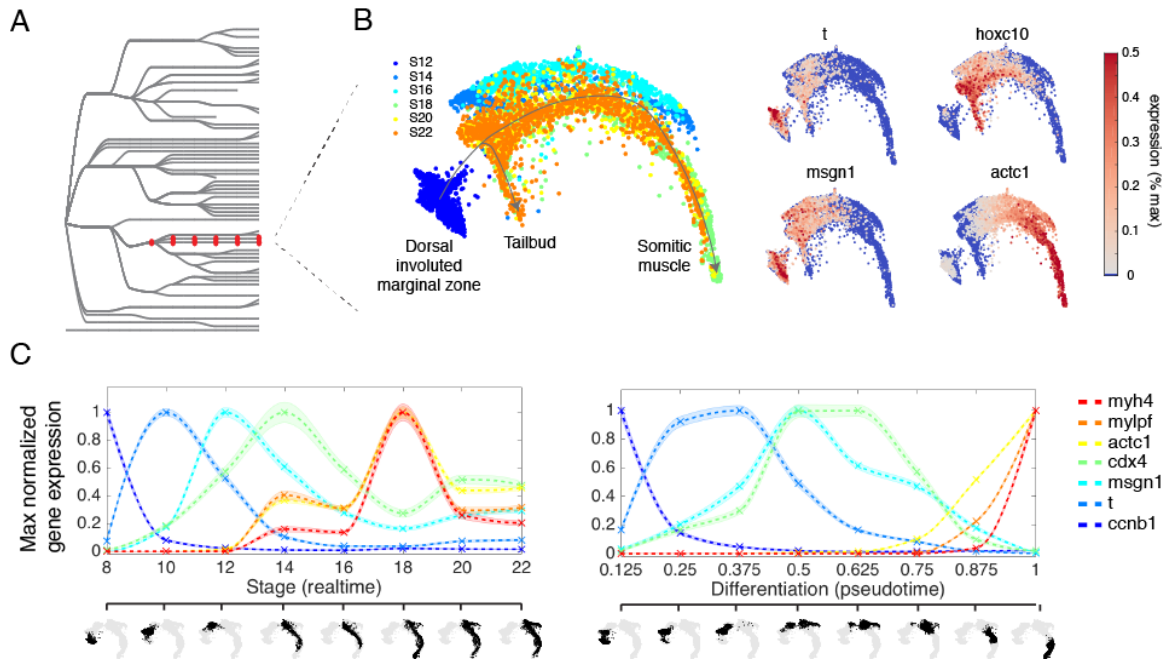


Fig. S10. Differentiation asynchrony in the somites.

(A) Cell state tree from Fig. 2C, colored to indicate the presomitic mesoderm and somitic muscle, which form parallel branches due to differentiation asynchrony. (B) Visualization of single cells (replicate 1) belonging to both branches reveals the underlying asynchronous differentiation process. Cells are colored either by timepoint (left), or by expression of key differentiation genes (right). Cells organize by differentiation progress with immature presomitic mesoderm states on the left (*t*, *msgn1*, *hoxc10*), which then differentiate into more mature somitic tissue (*actc1*). Mixing of early (S14) and late (S22) cells indicates asynchrony. (C) Real time versus pseudotime gene expression dynamics in somitic muscle differentiation. Asynchrony causes early, middle, and late waves of gene expression to blur compared to a pseudotemporal ordering, which respects a strictly ordered gene expression progression.

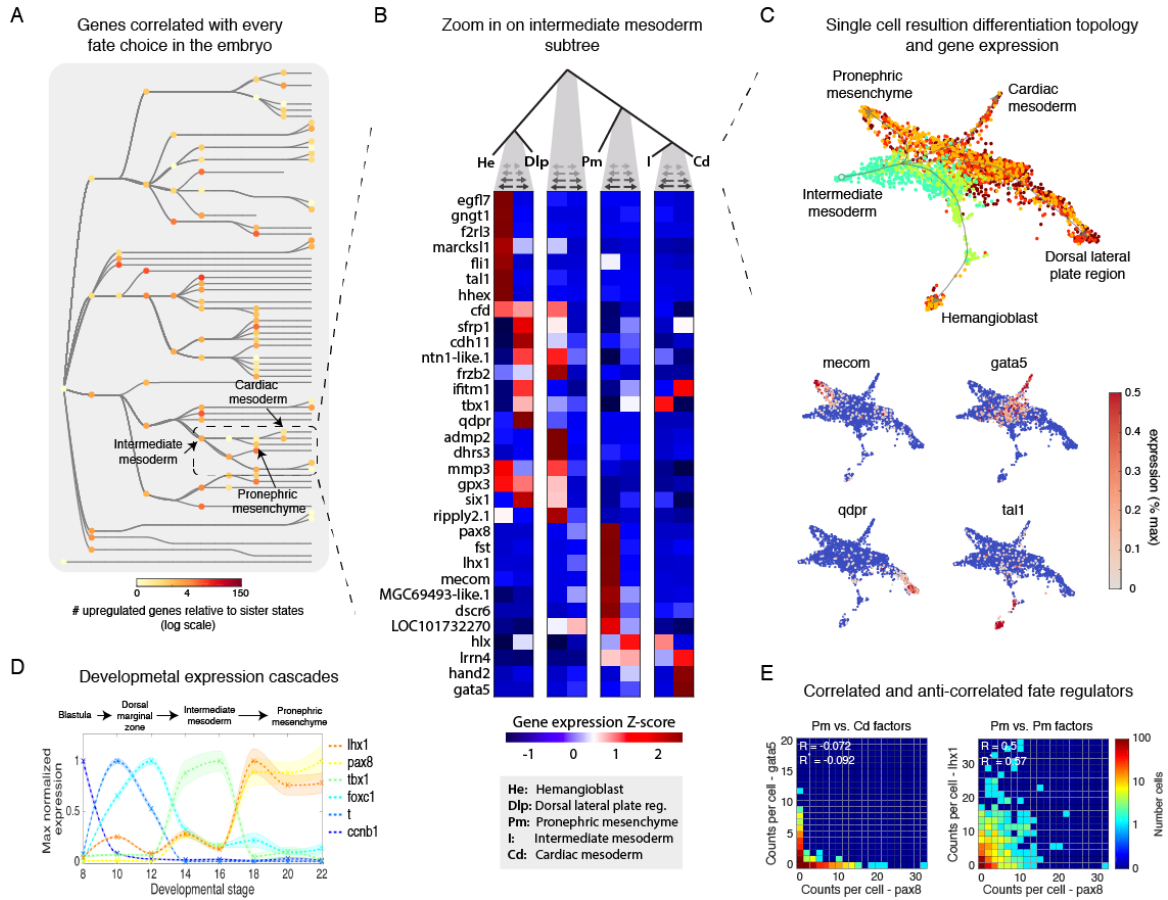


Fig. S11. The *Xenopus* cell state tree is a powerful gene expression resource.

(A) Global visualization of differential gene expression detected at every cell fate choice, showing the number of cell state specific genes (>4-fold enriched) relative to sister fates at each split. Tens to hundreds of fate-correlated genes are detected at every split. Panels B-D provide a more granular view of some of the gene expression information encoded within the intermediate mesoderm subtree (boxed). (B) Heatmap showing differentially expressed genes at each cell fate choice in the intermediate mesoderm. (C) Visualization of single cells (replicate 1) in the intermediate mesoderm subtree, showcasing examples of lineage specific TFs (also shown in Fig. 2E,F). (D) Reconstruction of the gene expression dynamics occurring during pronephric mesenchyme development. Cascades of early, middle, and late gene expression can be detected. (E) scRNA-seq data allows interrogation of gene coexpression in single cells, revealing correlated and anti-correlated gene pairs during fate choices. Pm, Cd as defined in grey box (panel B).

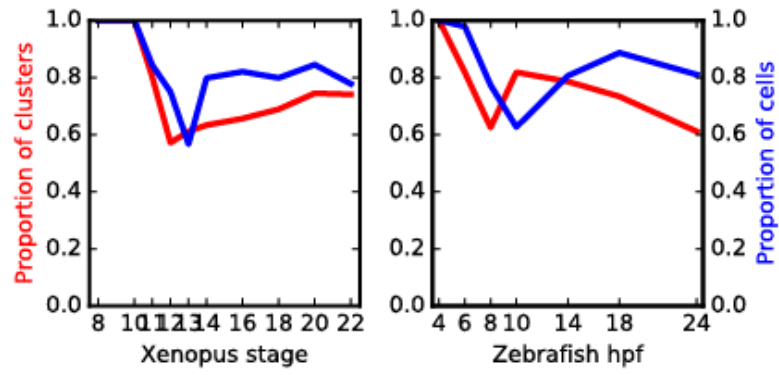


Fig. S12. Proportion of cell states and total cells matched between Xenopus and Zebrafish cell state trees.

Proportion of matched states (left) or matched cells (right) from the orthologous cell state mapping in Fig. 4A.

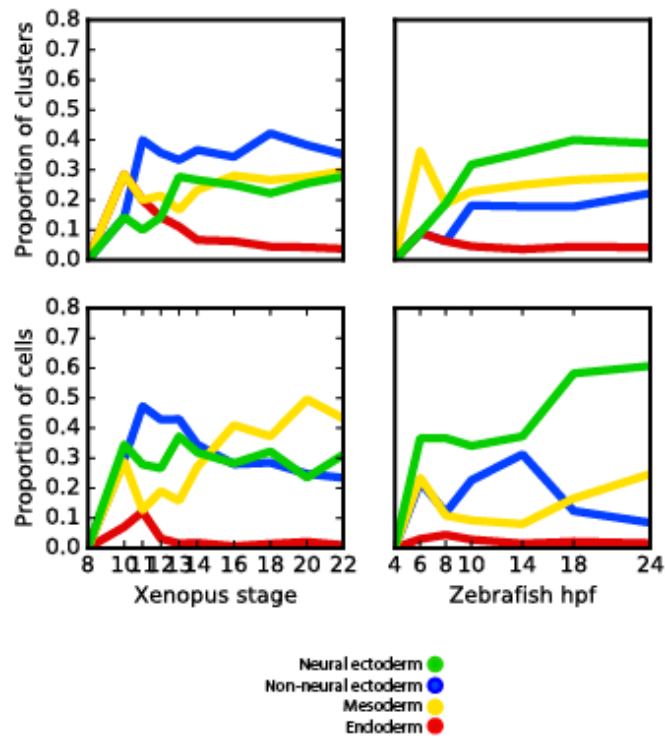


Fig. S13. Changes in germ layer proportions between Xenopus and Zebrafish.

Of sequenced cells, zebrafish are 60% neuroectoderm by 24hpf, as compared to 31% in stage 22 Xenopus embryos. Xenopus embryos by contrast have >2-fold more non-neural ectoderm cells (23% compared to 9%) at the same stage.

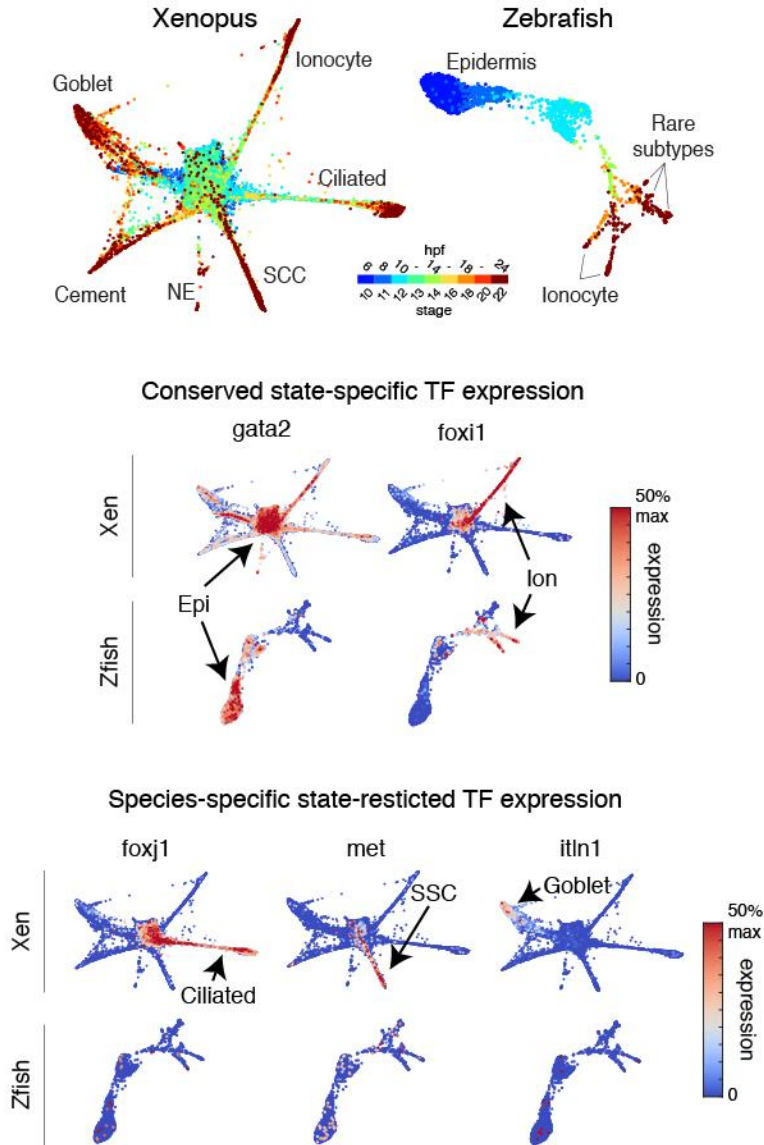


Fig. S14. Conserved and species specific gene expression on matched epidermal subtrees from Xenopus and Zebrafish.

This figure shows gene expression on the SPRING plots shown in Fig. 4B. The Xenopus and Zebrafish share conserved epidermal and ionocyte cell states, marked by conserved expression of *gata2* and *foxi1* respectively. By contrast, species specific states such as ciliated epidermal cells, small secretory cells (SSC), and goblet cells in Xenopus, are show species specific marker gene expression.

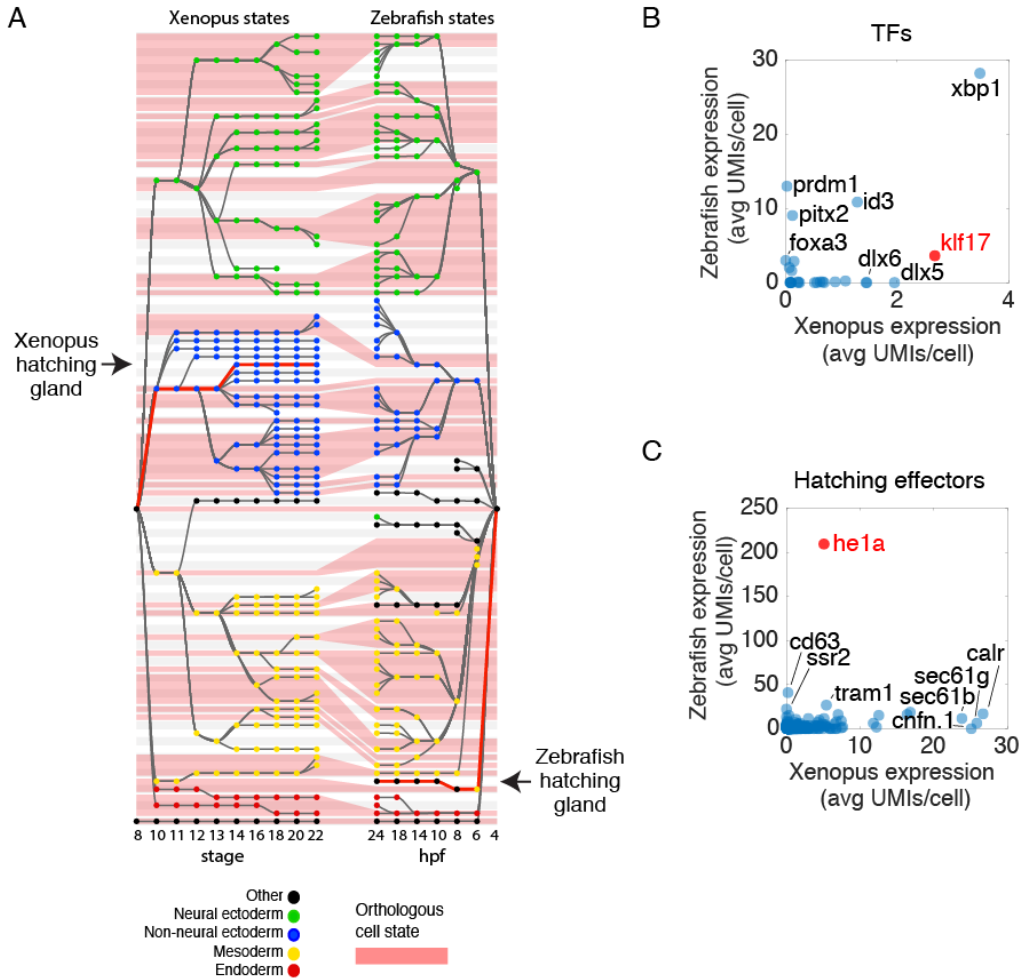
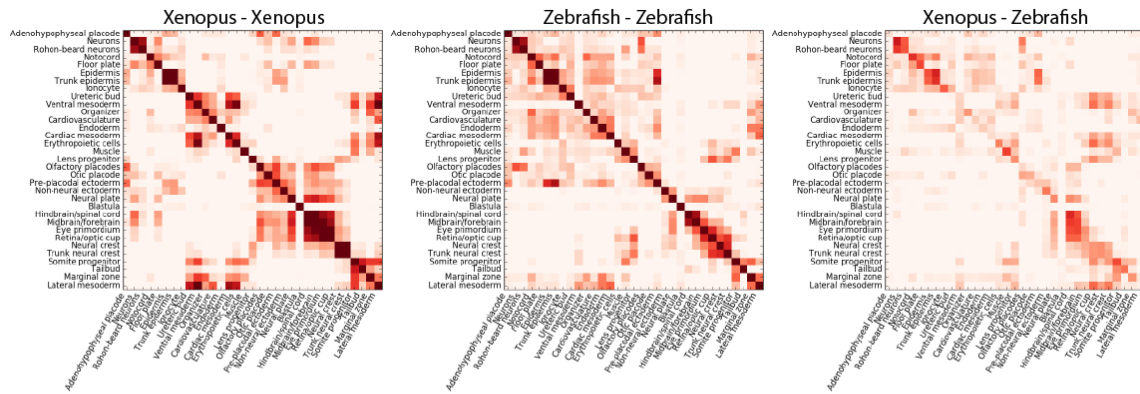


Fig. S15. scRNA-seq confirms alternative germ layer origins and conserved gene expression in the hatching gland of frog and fish.

(A) Aligned Xenopus and Zebrafish cell state trees (Fig. 4A), indicating the hatching gland differentiation path in both species (red). The unsupervised tree building algorithm correctly identifies that the hatching gland emerges from non-neural ectoderm in frog but mesodermal organizer tissue in fish. (B-C) Comparison of gene expression in the hatching glands of both species. scRNA-seq detects conserved expression of the TF, *klf17*, which is necessary for hatching gland differentiation in both species (B), as well as conserved expression of the hatching enzyme, *he1a* (C). Novel conserved and species specific gene expression is also detected in both species.

Correlation across genes with self-correlation > 0.5 in both species (N=3017)



Correlation across genes with self-correlation > 0.5 and inter-species correlation > 0.5 (N=770)

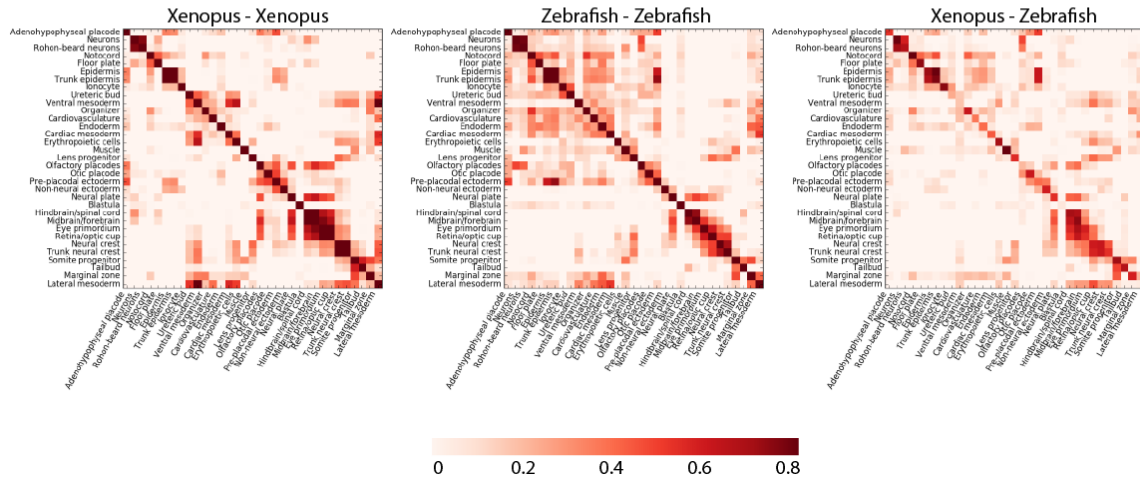


Fig. S16. Global gene expression similarity of matched Xenopus and Zebrafish cell states.

Heatmaps show the Pearson correlation between matched cell states (replicate 1) for the indicated species-comparison and pre-filtering stringency. Genes that are predictive of cell state within one species tend to be poorly predictive across species. Self-correlation is computed by jackknifing the data within a species (see Materials and Methods).

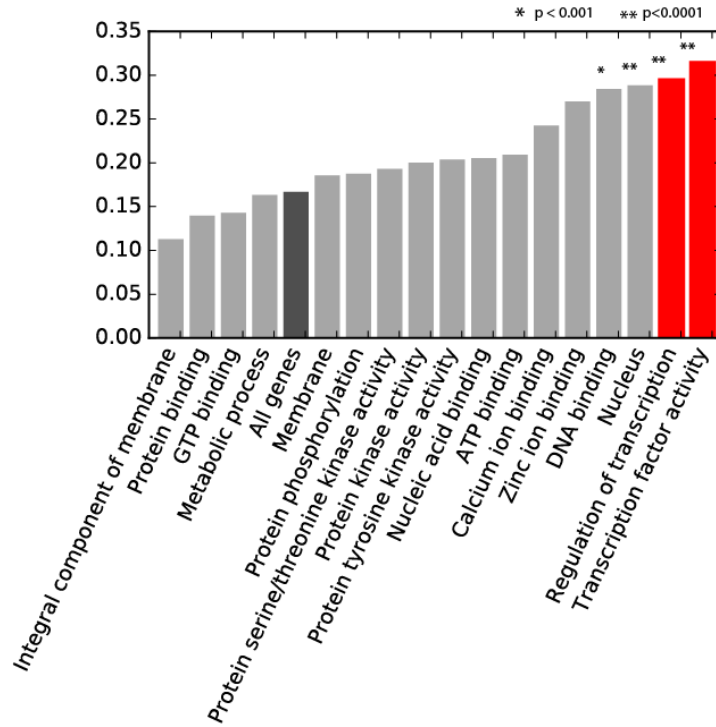


Fig. S17. MLP genes show enrichment of TF-associated GO terms compared to all DE genes.

MLP genes were compared to all DE genes (dark grey bar) across embryonic fate choices in *Xenopus*. Enriched TF-associated terms include ‘Transcription factor activity’, ‘Regulation of transcription’, ‘Nucleus’, and ‘DNA binding’. P-values indicate significance as determined by a binomial test with Bonferroni correction.

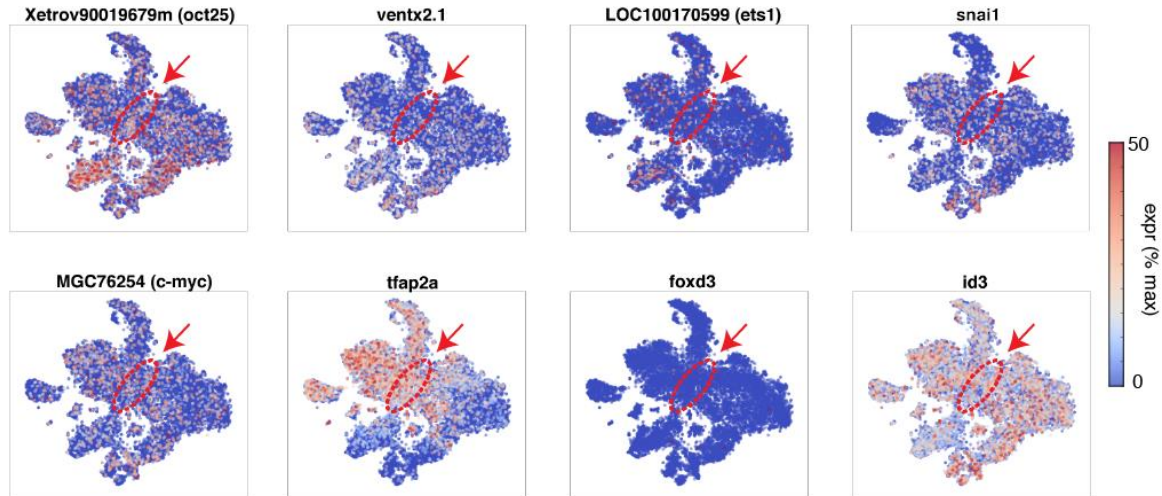


Fig. S18. Expression at stage 11 of individual genes proposed to constitute a retained blastula pluripotency circuit in neural crest cells.

This figure shows the individual genes comprising the ‘pluri-circuit’ score plotted in main text **Fig. 7E**. Each gene is expressed broadly, including in non-pluripotent tissues. The red circled region indicates the boundary between neuroectoderm and non-neural ectoderm cells, that associates with neural plate border by S12, and neural crest by S13. This region shows no special enrichment of these factors. Foxd3 is lowly expressed by scRNA-seq but detected in mesoderm and dorsal endoderm as well as ectoderm by in situ (Xenbase), showing that it is not a special case.

Table S1. scRNA-seq data summary.

Sample information for all scRNA-seq libraries. Libraries from replicate 2-3, which were sequenced to a lower average depth, are shaded in grey. # Libs indicates the number of separate droplet-emulsions per sample. Raw cells indicates the number of cells that passed the indicated min counts UMI threshold. Excluded cells were removed after identification as low count clusters, or as cell doublets (see methods section “Data clean up”). Final cells passed all filters and are incorporated into the cell state tree.

Sample	Biological replicate ID	# Libs	Min counts	Avg. UMIs / cell	Raw cells	Excluded cells (%)	Final cells
S8	1	2	1,000	3,661	875	0 (0%)	875
S8	2	1	500	813	318	0 (0%)	318
S10	1	3	1,000	4,517	6,507	580 (9%)	5,927
S11	2	6	400	890	6,078	272 (4%)	5,806
S11-NPB	3	4	400	756	9,348	512 (5%)	8,836
S12	1	3	1,000	4,713	6,091	713 (12%)	5,378
S12	2	3	400	881	7,722	1,282 (17%)	6,440
S13	2	4	400	896	8,931	1,096 (12%)	7,835
S14	1	4	1,000	5,948	4,910	451 (9%)	4,459
S14	2	3	300	913	5,885	1,854 (32%)	4,031
S16	1	4	1,000	5,023	5,173	302 (6%)	4,871
S16	2	3	500	1,412	8,305	1,096 (13%)	7,209
S18	1	4	1,000	5,839	5,130	251 (5%)	4,879
S18	2	3	500	1,415	7,302	785 (11%)	6,517
S20	1	5	1,000	5,744	6,036	515 (9%)	5,521
S20	2	3	500	1,590	10,606	853 (8%)	9,753
S22	1	6	1,000	5,674	8,113	541 (7%)	7,572
S22	2	14	500	1,875	29,636	2,873 (10%)	26,763
Total / avg:				2,920	136,966	13,976 (10%)	122,990

Movie S1. Collection of dissociated *Xenopus* embryo cells by inDrop microfluidics – cell inlet.

Slow motion recording of *Xenopus* cells entering the inDrop microfluidic device. Cells are round dark spheres, characteristic of good condition, and show significant variability in size as expected.

Movie S2. Collection of dissociated *Xenopus* embryo cells by inDrop microfluidics – droplet junction.

Slow motion recording of droplet-formation in the inDrop microfluidic device, showing co-encapsulation of *Xenopus* cells (dark spheres, entering from middle right) with barcoding hydrogels (larger transparent spheres, entering from bottom right) at ~4s and ~7s. Reverse transcription reagents enter from the top right channel and mix with the aqueous buffer containing cells. Oil flows from the top- and bottom-middle and pinches the two aqueous phases to form droplets. Each drop is ~3nL.

Movie S3. Collection of dissociated *Xenopus* embryo cells by inDrop microfluidics – collection outlet.

Slow motion recording of droplets entering the collection outlet. 3 cells are visibly co-encapsulated with hydrogel beads.

Additional Data Table S1 (separate file). Literature supported marker genes and XAO term associations for each *Xenopus* embryonic cell state.

Table with 2,930 rows by 4 columns, provided as .txt file with tab-delimited values. First row is column headers. Column 1: Cell state name matching the 87 annotations in Figure 2. Column 2: Best match Xenbase name for each cell state in column 1, referred to by XAO identifier, that is used to identify literature supported marker genes. Column 3: Marker gene name. Column 4: PubMed ID (PMID) for a published paper reporting that the gene in column 3 is a marker for the cell state in column 2. See Materials and Methods section “Annotation of cell states” for further details.

Additional Data Table S2 (separate file). Benchmarking of each edge on *Xenopus* cell state tree against XAO developmental lineage relationships.

Table with 258 rows by 9 columns (A-I), provided as a .txt file with tab separated values. First row is column headers. Overall the table assesses each of the 257 edges from state A (column A) to state B (column B) on the tree – shown in figure 2 of the main text – against the lineage relationships documented on Xenbase (referred to by XAO identifier) or the indicated paper (column C). The verdict for each edge is summarized in column D as either: ‘OK’, when the inferred edge agrees with Xenbase; ‘Discovery: early start’, when the generated state (column B) appears at least 2 stages earlier than reported on Xenbase; ‘Ends too late’, when the generated state (column B) exists at least 2 stages later than reported on Xenbase; ‘Ends too early’, when the ancestor state (column A) disappears at least 2 stages earlier than reported on Xenbase; or ‘Error in tree’, when the inferred edge is inconsistent with lineage relationships reported on Xenbase. Columns E and F record the start dates for each of the annotated states when they exist in both the tree in figure 2 and on Xenbase, respectively; these values were used to make Figure 3A of the main text. Columns H and I tally and summarize the validation results. See Materials and Methods section “Benchmarking of cell state tree similarity relationships against XAO” for further details.

Additional Data Table S3 (separate file). Differentially expressed genes at every cell fate choice identified in early *Xenopus* development.

Table with 1,514 rows by 5 columns, provided as .txt file with comma separated values. The table documents genes that are enriched in each state (column 1) compared to its sister states – i.e. those that have the same parent – at every fate choice, or splitting event, on the cell state tree shown in Figure 2. These genes are candidate regulators of the cell fate choice that promote differentiation into the indicated state (column 1). First row is column headers. Column 1: state name. Column 2: gene name. Column 3: fold-change compared to sister states (pooled if more than one). Column 4: average expression in state from column 1. Column 5: p-value for enriched expression. Only genes >2-fold enriched, average detected expression >1UMI/cell, and adjusted p-value <0.0001 are reported. See Materials and Methods section “Differential gene expression at cell fate choices” for further details.

Additional Data Table S4 (separate file). Automatically identified top marker genes for each embryonic *Xenopus* cell state.

Table with 2,159 rows by 6 columns, provided as .txt file with comma separated values. The table systematically documents genes that are highly-specific for each of the 72 cell state annotations shown in Figure 2. First row is column headers. Column 1: state name. Column 2: marker gene name. Column 3: local enrichment – i.e. average expression in state from column 1 divided by the average expression of the state with the next highest expression. Column 4: global enrichment – i.e. average expression in state from column 1 divided by the average expression in the rest of the embryo. Column 5: average expression in UMIs/cell for the state in column 1. Column 6: fraction of cells in state indicated in column one with detected expression of the marker gene. Marker genes for each state are ranked by their local enrichment as this is a good measure of the desired attributes of a marker gene for methods such as in situ staining. See Materials and Methods section “Identification of marker genes for each cell state” for further details on the criteria used to define marker genes.

Additional Data Table S5 (separate file). Catalog of re-used TFs and their tissue associations in early *Xenopus* development.

Table with 74 rows and 8 columns. Row 1 gives column headers. Column 1 specifies each of the 74 reused TFs identified during the analysis presented in Fig. 5. Columns 2-8 specify each tissue where an initiation event of each TF was identified.

Additional Data Table S6 (separate file). Catalog of multi-lineage primed gene pairs at fate choices in early *Xenopus* development.

Table with 1,386 rows and 5 columns. Overall the table describes 16 fate branch points with 412 MLP genes, in 1,385 MLP pairs. Row 1 gives column headers. Column 1 defines fate branch point where MLP genes were searched for. Column 2-3 give the names of 2 genes that show MLP at the fate choice noted in column 1. Column 4 gives the maximum coexpression of the gene pair in the subtree. Column 5 gives the minimum coexpression of the gene pair in the subtree, which occurs a timepoint after the maximum coexpression.

Additional Data Table S7 (separate file). Matching of X. Tr 9.0 transcripts to human gene symbols.

Table with 26,551 rows by 4 columns, provided as .txt file with tab-delimited values. Table documents matching between X. Tr 9.0 genome assembly gene symbols (column 1) with best match human gene symbols (column 2), as defined by a custom reciprocal protein blast pipeline (see Materials and Methods section “Expansion of gene symbol assignments in the XTr 9.0 reference transcriptome”). Columns 3 and 4 provide the forward and reverse blast e-values.

References

1. A. M. Klein, L. Mazutis, I. Akartuna, N. Tallapragada, A. Veres, V. Li, L. Peshkin, D. A. Weitz, M. W. Kirschner, Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell* **161**, 1187–1201 (2015). [doi:10.1016/j.cell.2015.04.044](https://doi.org/10.1016/j.cell.2015.04.044) [Medline](#)
2. J. Newport, M. Kirschner, A major developmental transition in early *Xenopus* embryos: I. characterization and timing of cellular changes at the midblastula stage. *Cell* **30**, 675–686 (1982). [doi:10.1016/0092-8674\(82\)90272-0](https://doi.org/10.1016/0092-8674(82)90272-0) [Medline](#)
3. J. Newport, M. Kirschner, A major developmental transition in early *Xenopus* embryos: II. Control of the onset of transcription. *Cell* **30**, 687–696 (1982). [doi:10.1016/0092-8674\(82\)90273-2](https://doi.org/10.1016/0092-8674(82)90273-2) [Medline](#)
4. J. Heasman, Patterning the early *Xenopus* embryo. *Development* **133**, 1205–1217 (2006). [doi:10.1242/dev.02304](https://doi.org/10.1242/dev.02304) [Medline](#)
5. K. Karimi, J. D. Fortriede, V. S. Lotay, K. A. Burns, D. Z. Wang, M. E. Fisher, T. J. Pells, C. James-Zorn, Y. Wang, V. G. Ponferrada, S. Chu, P. Chaturvedi, A. M. Zorn, P. D. Vize, Xenbase: A genomic, epigenomic and transcriptomic model organism database. *Nucleic Acids Res.* **46**, D861–D868 (2018). [Medline](#)
6. R. L. Davis, M. W. Kirschner, The fate of cells in the tailbud of *Xenopus laevis*. *Development* **127**, 255–267 (2000). [Medline](#)
7. D. E. Wagner, C. Weinreb, Z. M. Collins, J. A. Briggs, S. G. Megason, A. M. Klein, Single-cell mapping of gene expression landscapes and lineage in the zebrafish embryo. *Science* **360**, aar4362 (2018).
8. E. Segerdell, V. G. Ponferrada, C. James-Zorn, K. A. Burns, J. D. Fortriede, W. M. Dahdul, P. D. Vize, A. M. Zorn, Enhanced XAO: The ontology of *Xenopus* anatomy and development underpins more accurate annotation of gene expression and queries on Xenbase. *J. Biomed. Semantics* **4**, 31 (2013). [doi:10.1186/2041-1480-4-31](https://doi.org/10.1186/2041-1480-4-31) [Medline](#)
9. E. Segerdell, J. B. Bowes, N. Pollet, P. D. Vize, An ontology for *Xenopus* anatomy and development. *BMC Dev. Biol.* **8**, 92 (2008). [doi:10.1186/1471-213X-8-92](https://doi.org/10.1186/1471-213X-8-92) [Medline](#)
10. C. Trapnell, D. Cacchiarelli, J. Grimsby, P. Pokharel, S. Li, M. Morse, N. J. Lennon, K. J. Livak, T. S. Mikkelsen, J. L. Rinn, The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014). [doi:10.1038/nbt.2859](https://doi.org/10.1038/nbt.2859) [Medline](#)
11. C. Weinreb, S. Wolock, B. K. Tusi, M. Socolovsky, A. M. Klein, Fundamental limits on dynamic inference from single-cell snapshots. *Proc. Natl. Acad. Sci. U.S.A.* **115**, E2467–E2476 (2018). [doi:10.1073/pnas.1714723115](https://doi.org/10.1073/pnas.1714723115) [Medline](#)
12. G. Schiebinger *et al.*, Reconstruction of developmental landscapes by optimal-transport analysis of single-cell gene expression sheds light on cellular reprogramming. bioRxiv 191056 [preprint]. 27 September 2017. <https://doi.org/10.1101/191056>
13. M. Walmsley, A. Ciau-Uitz, R. Patient, Adult and embryonic blood and endothelium derive from distinct precursor populations which are differentially programmed by BMP in *Xenopus*. *Development* **129**, 5683–5695 (2002). [doi:10.1242/dev.00169](https://doi.org/10.1242/dev.00169) [Medline](#)

14. F. Liu, M. Walmsley, A. Rodaway, R. Patient, Fli1 acts at the top of the transcriptional network driving blood and endothelial development. *Curr. Biol.* **18**, 1234–1240 (2008). [doi:10.1016/j.cub.2008.07.048](https://doi.org/10.1016/j.cub.2008.07.048) [Medline](#)
15. T. Nagasawa, M. Kawaguchi, T. Yano, K. Sano, M. Okabe, S. Yasumasu, Evolutionary Changes in the Developmental Origin of Hatching Gland Cells in Basal Ray-Finned Fishes. *Zool. Sci.* **33**, 272–281 (2016). [doi:10.2108/zs150183](https://doi.org/10.2108/zs150183) [Medline](#)
16. M. Yasutomi, T. Hama, Electron microscopic study on the xanthophore differentiation in *Xenopus laevis*, with special reference to their pterinosomes. *J. Ultrastruct. Res.* **38**, 421–432 (1972). [doi:10.1016/0022-5320\(72\)90080-9](https://doi.org/10.1016/0022-5320(72)90080-9) [Medline](#)
17. J. A. Briggs, V. C. Li, S. Lee, C. J. Woolf, A. Klein, M. W. Kirschner, Mouse embryonic stem cells can differentiate via multiple paths to the same state. *eLife* **6**, e26945 (2017). [doi:10.7554/eLife.26945](https://doi.org/10.7554/eLife.26945) [Medline](#)
18. H. Weber, C. E. Symes, M. E. Walmsley, A. R. Rodaway, R. K. Patient, A role for GATA5 in *Xenopus* endoderm specification. *Development* **127**, 4345–4360 (2000). [Medline](#)
19. J. F. Reiter, J. Alexander, A. Rodaway, D. Yelon, R. Patient, N. Holder, D. Y. R. Stainier, Gata5 is required for the development of the heart and endoderm in zebrafish. *Genes Dev.* **13**, 2983–2995 (1999). [doi:10.1101/gad.13.22.2983](https://doi.org/10.1101/gad.13.22.2983) [Medline](#)
20. H. Weintraub, The MyoD family and myogenesis: Redundancy, networks, and thresholds. *Cell* **75**, 1241–1244 (1993). [doi:10.1016/0092-8674\(93\)90610-3](https://doi.org/10.1016/0092-8674(93)90610-3) [Medline](#)
21. T. Ravasi, H. Suzuki, C. V. Cannistraci, S. Katayama, V. B. Bajic, K. Tan, A. Akalin, S. Schmeier, M. Kanamori-Katayama, N. Bertin, P. Carninci, C. O. Daub, A. R. R. Forrest, J. Gough, S. Grimmond, J.-H. Han, T. Hashimoto, W. Hide, O. Hofmann, A. Kamburov, M. Kaur, H. Kawaji, A. Kubosaki, T. Lassmann, E. van Nimwegen, C. R. MacPherson, C. Ogawa, A. Radovanovic, A. Schwartz, R. D. Teasdale, J. Tegnér, B. Lenhard, S. A. Teichmann, T. Arakawa, N. Ninomiya, K. Murakami, M. Tagami, S. Fukuda, K. Imamura, C. Kai, R. Ishihara, Y. Kitazume, J. Kawai, D. A. Hume, T. Ideker, Y. Hayashizaki, An atlas of combinatorial transcriptional regulation in mouse and man. *Cell* **140**, 744–752 (2010). [doi:10.1016/j.cell.2010.01.044](https://doi.org/10.1016/j.cell.2010.01.044) [Medline](#)
22. J. Monod, F. Jacob, Teleonomic mechanisms in cellular metabolism, growth, and differentiation. *Cold Spring Harb. Symp. Quant. Biol.* **26**, 389–401 (1961). [doi:10.1101/SQB.1961.026.01.048](https://doi.org/10.1101/SQB.1961.026.01.048) [Medline](#)
23. D. E. Cohen, D. Melton, Turning straw into gold: Directing cell fate for regenerative medicine. *Nat. Rev. Genet.* **12**, 243–252 (2011). [doi:10.1038/nrg2938](https://doi.org/10.1038/nrg2938) [Medline](#)
24. M. Hu, D. Krause, M. Greaves, S. Sharkis, M. Dexter, C. Heyworth, T. Enver, Multilineage gene expression precedes commitment in the hemopoietic system. *Genes Dev.* **11**, 774–785 (1997). [doi:10.1101/gad.11.6.774](https://doi.org/10.1101/gad.11.6.774) [Medline](#)
25. M. Thomson, S. J. Liu, L.-N. Zou, Z. Smith, A. Meissner, S. Ramanathan, Pluripotency factors in embryonic stem cells regulate differentiation into germ layers. *Cell* **145**, 875–889 (2011). [doi:10.1016/j.cell.2011.05.017](https://doi.org/10.1016/j.cell.2011.05.017) [Medline](#)
26. L. Velten, S. F. Haas, S. Raffel, S. Blaszkiewicz, S. Islam, B. P. Hennig, C. Hirche, C. Lutz, E. C. Buss, D. Nowak, T. Boch, W.-K. Hofmann, A. D. Ho, W. Huber, A. Trumpp, M. A.

- G. Essers, L. M. Steinmetz, Human haematopoietic stem cell lineage commitment is a continuous process. *Nat. Cell Biol.* **19**, 271–281 (2017). [doi:10.1038/ncb3493](https://doi.org/10.1038/ncb3493) [Medline](#)
27. J. Shu, C. Wu, Y. Wu, Z. Li, S. Shao, W. Zhao, X. Tang, H. Yang, L. Shen, X. Zuo, W. Yang, Y. Shi, X. Chi, H. Zhang, G. Gao, Y. Shu, K. Yuan, W. He, C. Tang, Y. Zhao, H. Deng, Induction of pluripotency in mouse somatic cells with lineage specifiers. *Cell* **153**, 963–975 (2013). [doi:10.1016/j.cell.2013.05.001](https://doi.org/10.1016/j.cell.2013.05.001) [Medline](#)
28. P. Laslo, C. J. Spooner, A. Warmflash, D. W. Lancki, H.-J. Lee, R. Sciammas, B. N. Gantner, A. R. Dinner, H. Singh, Multilineage transcriptional priming and determination of alternate hematopoietic cell fates. *Cell* **126**, 755–766 (2006). [doi:10.1016/j.cell.2006.06.052](https://doi.org/10.1016/j.cell.2006.06.052) [Medline](#)
29. F. C. Wardle, J. C. Smith, Refinement of gene expression patterns in the early *Xenopus* embryo. *Development* **131**, 4687–4696 (2004). [doi:10.1242/dev.01340](https://doi.org/10.1242/dev.01340) [Medline](#)
30. E. Buitrago-Delgado, K. Nordin, A. Rao, L. Geary, C. LaBonne, Shared regulatory programs suggest retention of blastula-stage potential in neural crest cells. *Science* **348**, 1332–1335 (2015). [doi:10.1126/science.aaa3655](https://doi.org/10.1126/science.aaa3655) [Medline](#)
31. E. Lubeck, A. F. Coskun, T. Zhiyentayev, M. Ahmad, L. Cai, Single-cell in situ RNA profiling by sequential hybridization. *Nat. Methods* **11**, 360–361 (2014). [doi:10.1038/nmeth.2892](https://doi.org/10.1038/nmeth.2892) [Medline](#)
32. K. H. Chen, A. N. Boettiger, J. R. Moffitt, S. Wang, X. Zhuang, Spatially resolved, highly multiplexed RNA profiling in single cells. *Science* **348**, aaa6090 (2015). [doi:10.1126/science.aaa6090](https://doi.org/10.1126/science.aaa6090) [Medline](#)
33. E. Z. Macosko, A. Basu, R. Satija, J. Nemesh, K. Shekhar, M. Goldman, I. Tirosh, A. R. Bialas, N. Kamitaki, E. M. Martersteck, J. J. Trombetta, D. A. Weitz, J. R. Sanes, A. K. Shalek, A. Regev, S. A. McCarroll, Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202–1214 (2015). [doi:10.1016/j.cell.2015.05.002](https://doi.org/10.1016/j.cell.2015.05.002) [Medline](#)
34. N. Karaiskos, P. Wahle, J. Alles, A. Boltengagen, S. Ayoub, C. Kipar, C. Kocks, N. Rajewsky, R. P. Zinzen, The *Drosophila* embryo at single-cell transcriptome resolution. *Science* **358**, 194–199 (2017). [doi:10.1126/science.aan3235](https://doi.org/10.1126/science.aan3235) [Medline](#)
35. J. Cao, J. S. Packer, V. Ramani, D. A. Cusanovich, C. Huynh, R. Daza, X. Qiu, C. Lee, S. N. Furlan, F. J. Steemers, A. Adey, R. H. Waterston, C. Trapnell, J. Shendure, Comprehensive single-cell transcriptional profiling of a multicellular organism. *Science* **357**, 661–667 (2017). [doi:10.1126/science.aam8940](https://doi.org/10.1126/science.aam8940) [Medline](#)
36. K. Achim, J.-B. Pettit, L. R. Saraiva, D. Gavriouchkina, T. Larsson, D. Arendt, J. C. Marioni, High-throughput spatial mapping of single-cell RNA-seq data to tissue of origin. *Nat. Biotechnol.* **33**, 503–509 (2015). [doi:10.1038/nbt.3209](https://doi.org/10.1038/nbt.3209) [Medline](#)
37. R. Satija, J. A. Farrell, D. Gennert, A. F. Schier, A. Regev, Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**, 495–502 (2015). [doi:10.1038/nbt.3192](https://doi.org/10.1038/nbt.3192) [Medline](#)
38. D. Forman, J. M. Slack, Determination and cellular commitment in the embryonic amphibian mesoderm. *Nature* **286**, 492–494 (1980). [doi:10.1038/286492a0](https://doi.org/10.1038/286492a0) [Medline](#)

39. J. Heasman, C. C. Wylie, P. Hausen, J. C. Smith, Fates and states of determination of single vegetal pole blastomeres of *X. laevis*. *Cell* **37**, 185–194 (1984). [doi:10.1016/0092-8674\(84\)90314-3](https://doi.org/10.1016/0092-8674(84)90314-3) [Medline](#)
40. S. Hontelez, I. van Kruijsbergen, G. Georgiou, S. J. van Heeringen, O. Bogdanovic, R. Lister, G. J. C. Veenstra, Embryonic transcription is controlled by maternally defined chromatin state. *Nat. Commun.* **6**, 10148 (2015). [doi:10.1038/ncomms10148](https://doi.org/10.1038/ncomms10148) [Medline](#)
41. R. C. Akkers, S. J. van Heeringen, U. G. Jacobi, E. M. Janssen-Megens, K.-J. François, H. G. Stunnenberg, G. J. C. Veenstra, A hierarchy of H3K4me3 and H3K27me3 acquisition in spatial gene regulation in *Xenopus* embryos. *Dev. Cell* **17**, 425–434 (2009). [doi:10.1016/j.devcel.2009.08.005](https://doi.org/10.1016/j.devcel.2009.08.005) [Medline](#)
42. M. Hemberger, W. Dean, W. Reik, Epigenetic dynamics of stem cells and cell lineage commitment: Digging Waddington's canal. *Nat. Rev. Mol. Cell Biol.* **10**, 526–537 (2009). [doi:10.1038/nrm2727](https://doi.org/10.1038/nrm2727) [Medline](#)
43. S. John, P. J. Sabo, R. E. Thurman, M.-H. Sung, S. C. Biddie, T. A. Johnson, G. L. Hager, J. A. Stamatoyannopoulos, Chromatin accessibility pre-determines glucocorticoid receptor binding patterns. *Nat. Genet.* **43**, 264–268 (2011). [doi:10.1038/ng.759](https://doi.org/10.1038/ng.759) [Medline](#)
44. L. Ho, G. R. Crabtree, Chromatin remodelling during development. *Nature* **463**, 474–484 (2010). [doi:10.1038/nature08911](https://doi.org/10.1038/nature08911) [Medline](#)
45. D. Arendt, J. M. Musser, C. V. H. Baker, A. Bergman, C. Cepko, D. H. Erwin, M. Pavlicev, G. Schlosser, S. Widder, M. D. Laubichler, G. P. Wagner, The origin and evolution of cell types. *Nat. Rev. Genet.* **17**, 744–757 (2016). [doi:10.1038/nrg.2016.127](https://doi.org/10.1038/nrg.2016.127) [Medline](#)
46. H. L. Sive, R. M. Grainger, R. M. Harland, *Early Development of Xenopus laevis: A Laboratory Manual* (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, 2000).
47. V. Savova, E. J. Pearl, E. Boke, A. Nag, I. Adzhubei, M. E. Horb, L. Peshkin, Transcriptomic insights into genetic diversity of protein-coding genes in *X. laevis*. *Dev. Biol.* **424**, 181–188 (2017). [doi:10.1016/j.ydbio.2017.02.019](https://doi.org/10.1016/j.ydbio.2017.02.019) [Medline](#)
48. UniProt Consortium, Update on activities at the Universal Protein Resource (UniProt) in 2013. *Nucleic Acids Res.* **41**, D43–D47 (2013). [Medline](#)
49. L. van der Maaten, G. Hinton, Visualizing Data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).
50. C. Weinreb, S. Wolock, A. M. Klein, SPRING: A kinetic interface for visualizing high dimensional single-cell expression data. *Bioinformatics* **34**, 1246–1248 (2018). [doi:10.1093/bioinformatics/btx792](https://doi.org/10.1093/bioinformatics/btx792) [Medline](#)
51. I. L. Blitz, K. D. Paraiso, I. Patrushev, W. T. Y. Chiu, K. W. Y. Cho, M. J. Gilchrist, A catalog of *Xenopus tropicalis* transcription factors and their regional expression in the early gastrula stage embryo. *Dev. Biol.* **426**, 409–417 (2017). [doi:10.1016/j.ydbio.2016.07.002](https://doi.org/10.1016/j.ydbio.2016.07.002) [Medline](#)