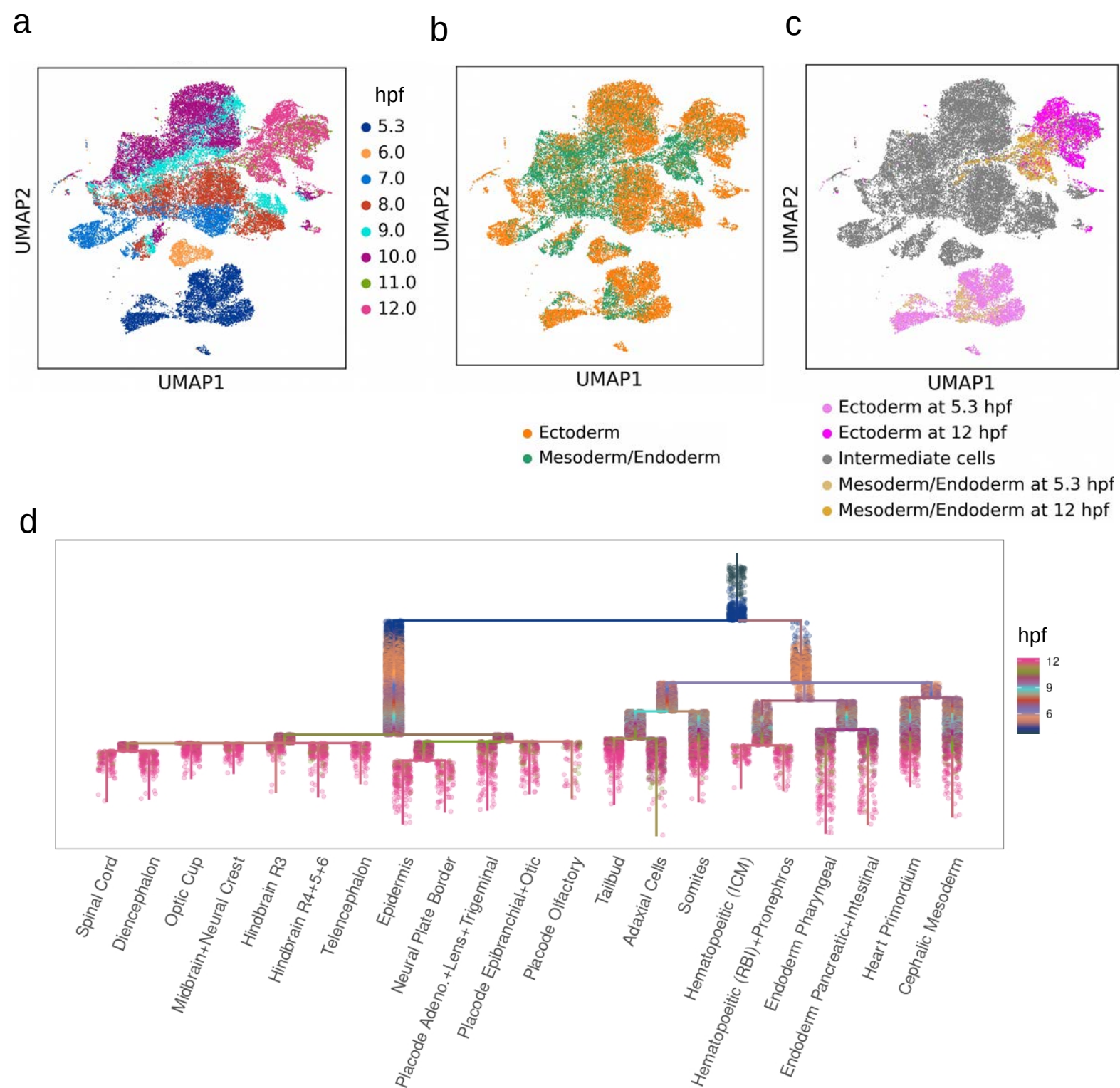


Cell Systems, Volume 15

Supplemental information

**Inference of differentiation trajectories
by transfer learning across biological processes**

Gaurav Jumde, Bastiaan Spanjaard, and Jan Philipp Junker



Supplementary Figure 1 | Developmental zebrafish dataset (related to Fig. 1).

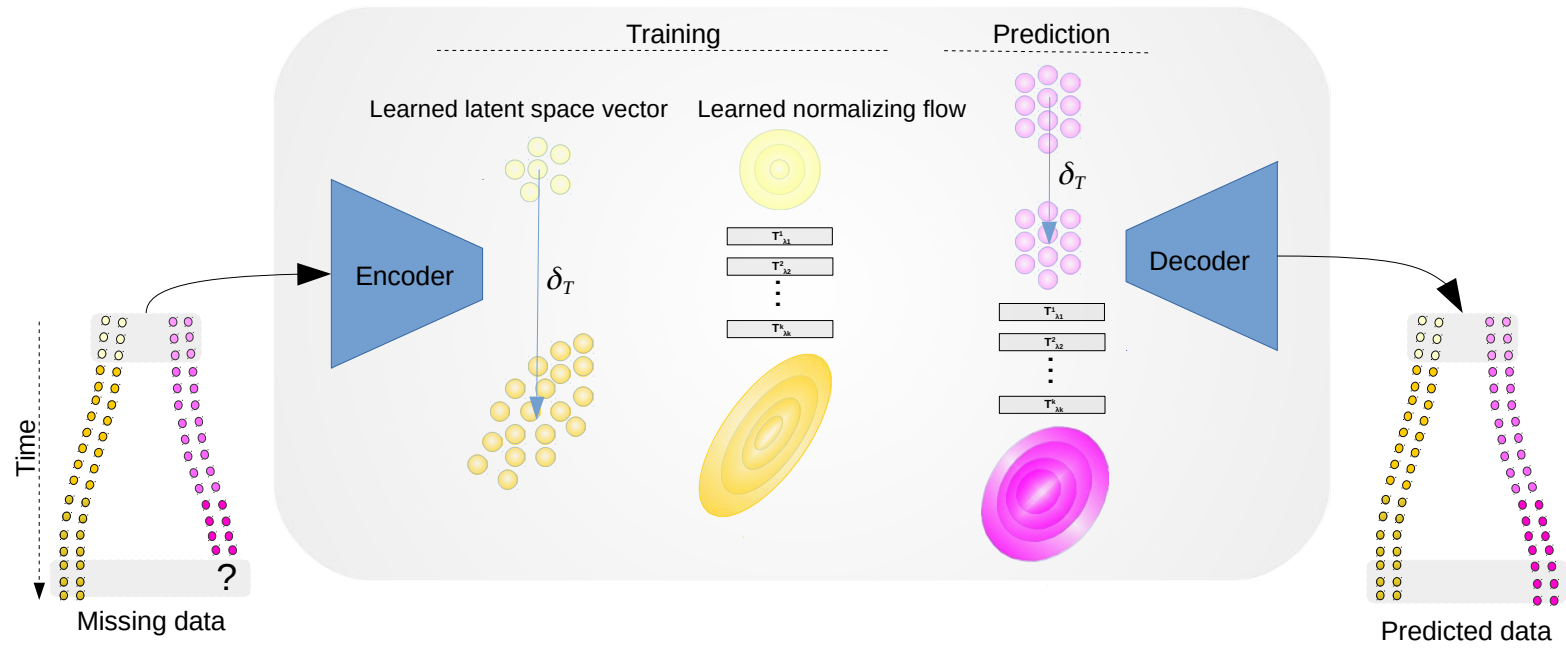
a-c) UMAP plot of the developmental zebrafish dataset, indicating hours post fertilization (hpf) **(a)**, germ layer **(b)**, and progenitor and mature cells at 5.3hpf and 12hpf of Ectoderm and Mesoderm/Endoderm germ layers **(c)**.

d) URD of the same dataset, indicating transcriptome-based inferred lineage splits. URD is the diffusion-based computational trajectory reconstruction method that was used in the original publication to analyze this dataset. Reproduced from original publication⁷.

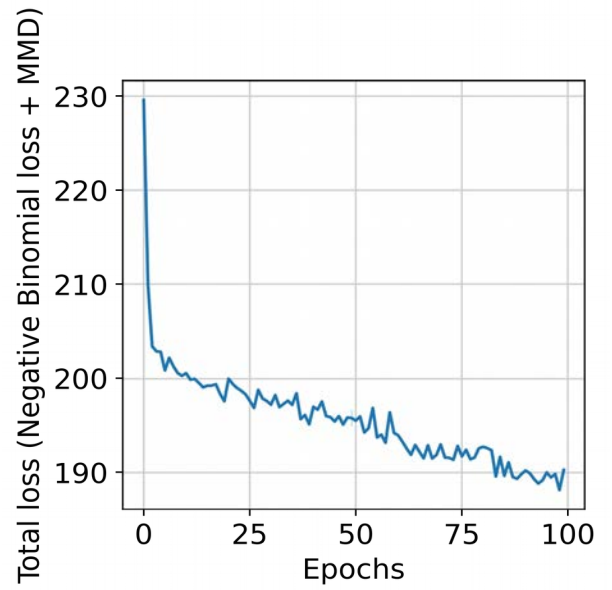
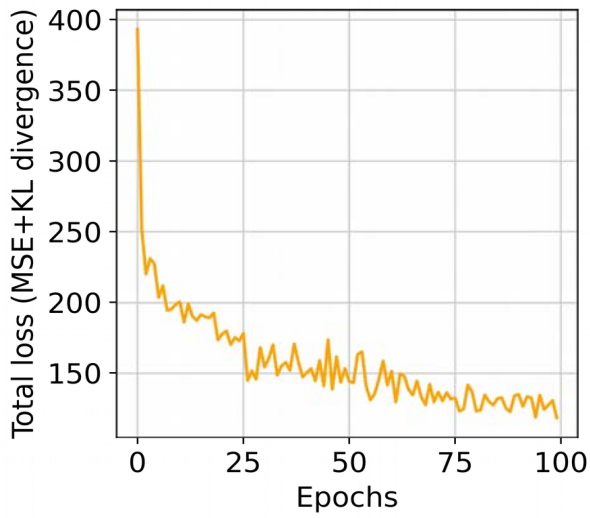
Data from Farrell et al., 2018⁷, GEO accession GSE106587.

DeepCellPredictor

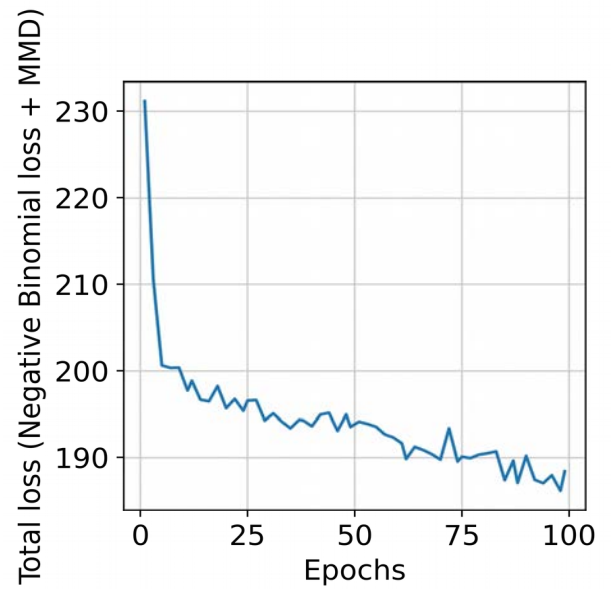
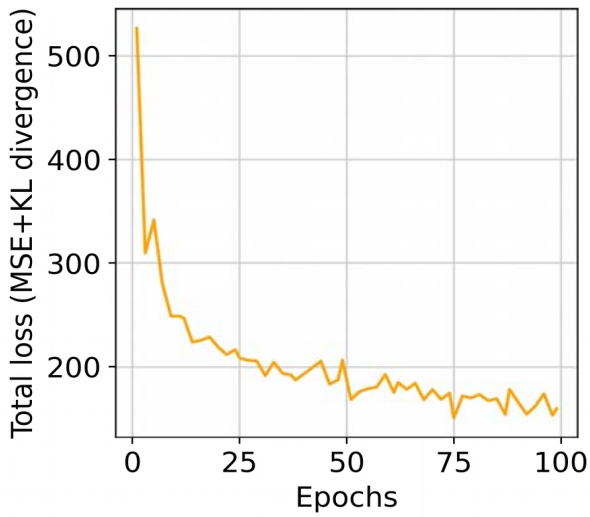
a



b



c



Supplementary Figure 2 | Training convergence of neural network models (related to Fig. 1).

a) Schematic of the DeepCellPredictor (DCP) model.

b-c) Total loss over epochs for regular variational autoencoder (left) and DCP model (right) trained on mesoderm/endoderm **(b)** and ectoderm data **(c)**.

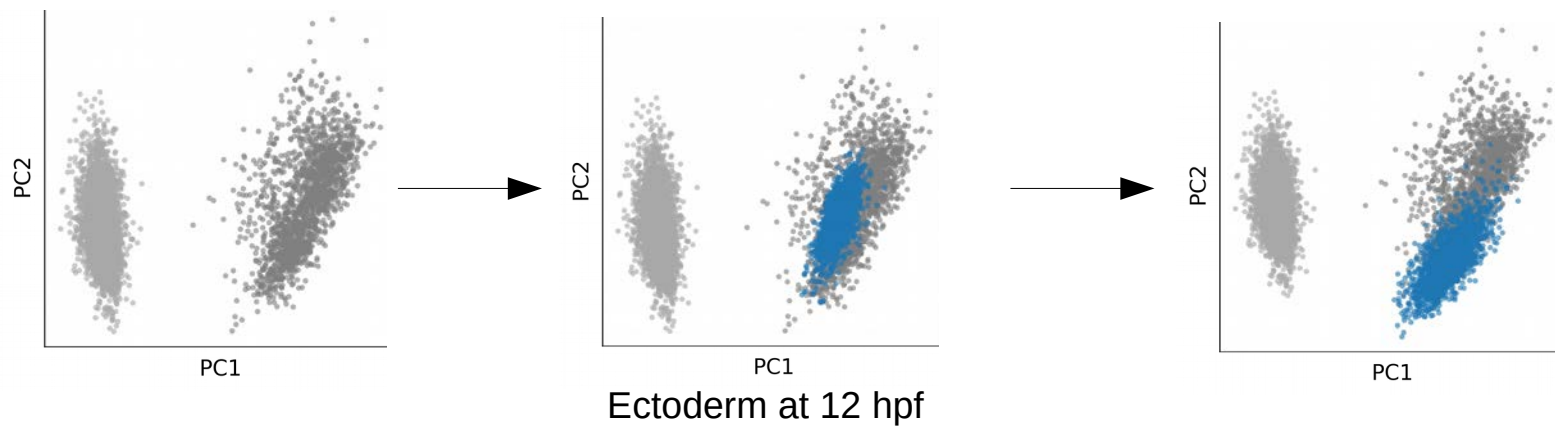
Data from Farrell et al., 2018⁷, GEO accession GSE106587.

Input data

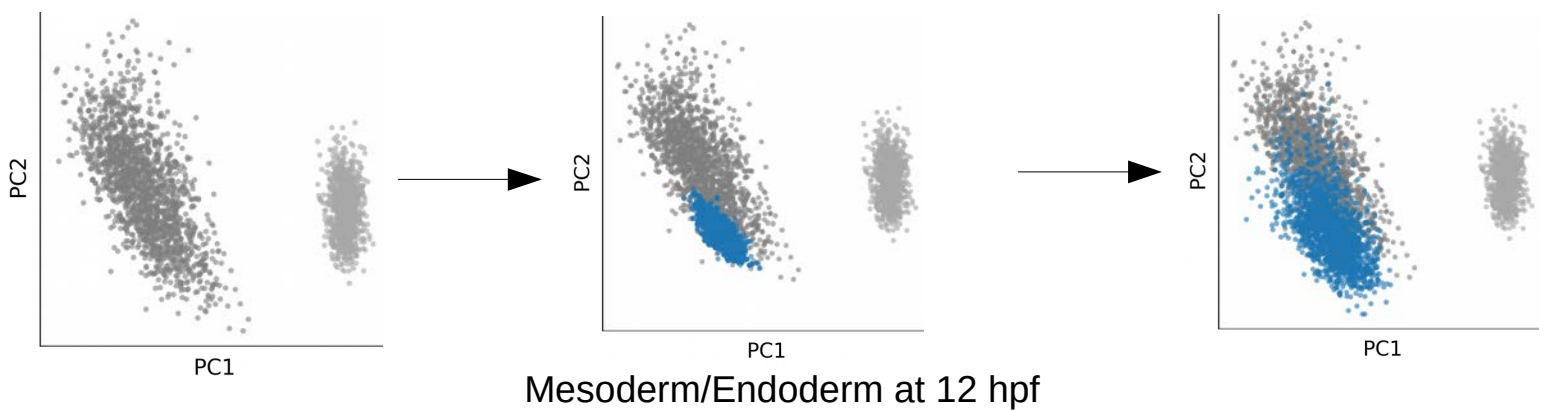
After vector arithmetic

After planar flows

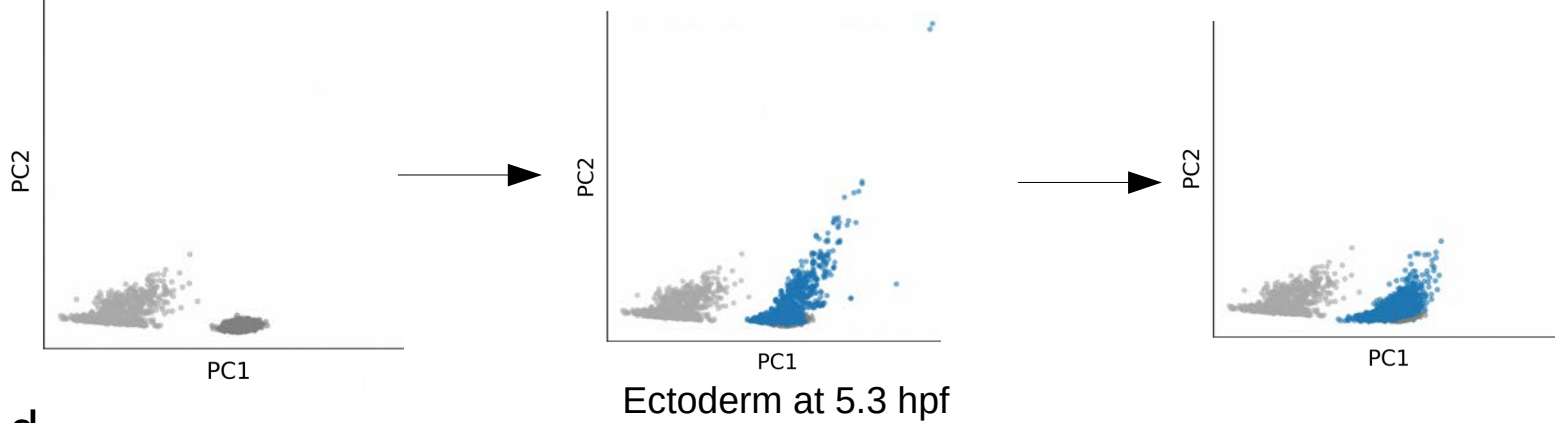
a



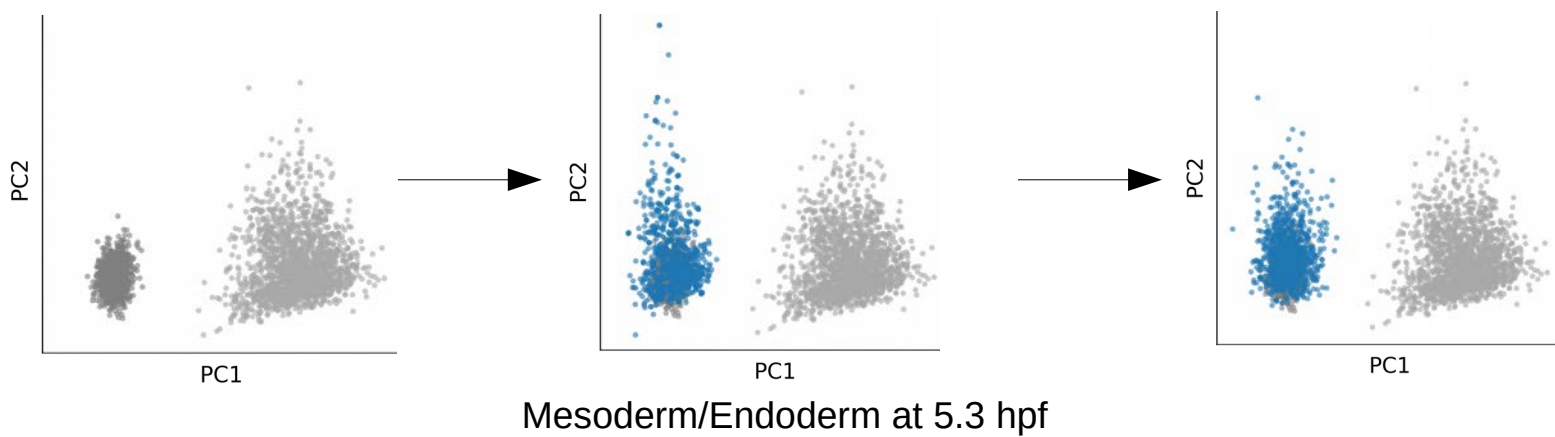
b



c



d

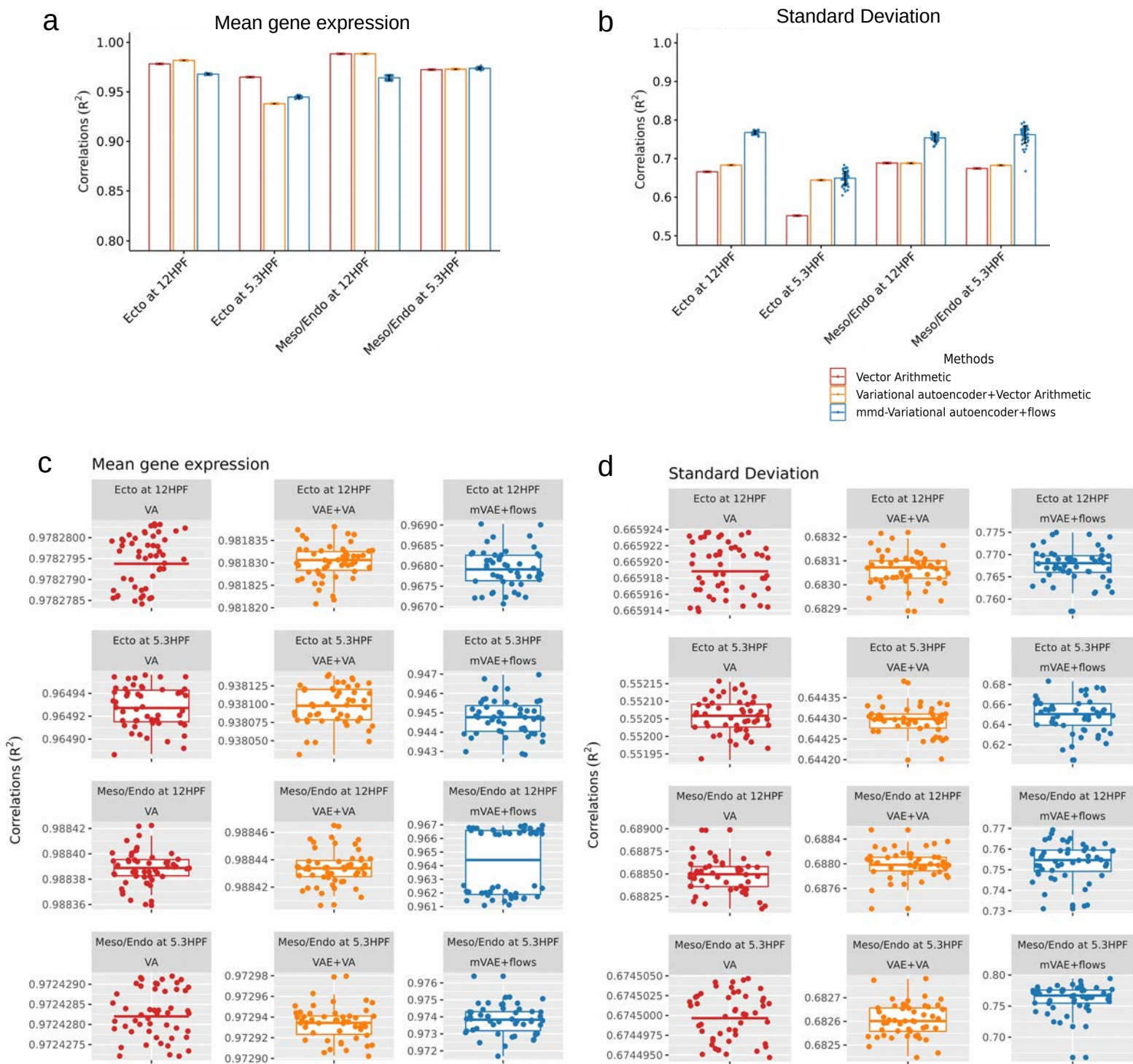


● Test data ● Real data ● Predicted data

Supplementary Figure 3 | Disentangling the effects of vector arithmetic and normalizing flows in DCP (related to Fig. 1).

To visualize the information transferred, we decoded latent space predictions after vector arithmetic and after normalizing flows and performed PCA for visualization. The predicted distribution of cells already starts resembling the target distribution after vector arithmetic due to our regularization approach, in contrast to the vector arithmetic performed in scGen (Fig. 1c). The normalizing flows allow us to further approximate the target distribution.

Data from Farrell et al., 2018⁷, GEO accession GSE106587.



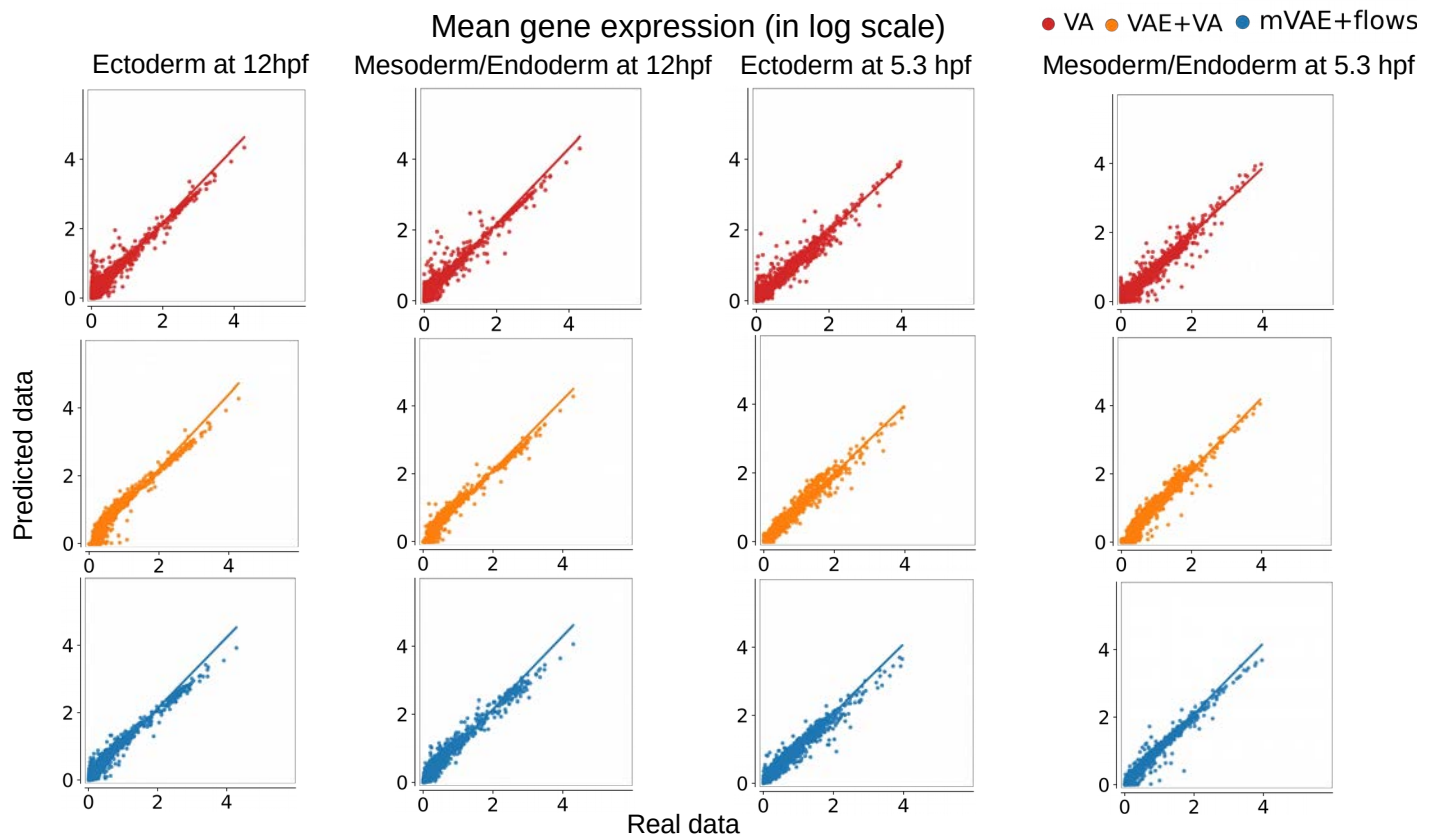
Supplementary Figure 4 | Comparison of prediction algorithms (vector arithmetic (VA), variational autoencoders with vector arithmetic (VAE+VA), and mmd-variational autoencoders with flows (mVAE+flows)) in the developing zebrafish (related to Fig. 2).

a-b) Correlation between real and predicted mean expression and expression variability. Error bars were determined by considering one standard deviation of uncertainty from the mean after 50 repeats.

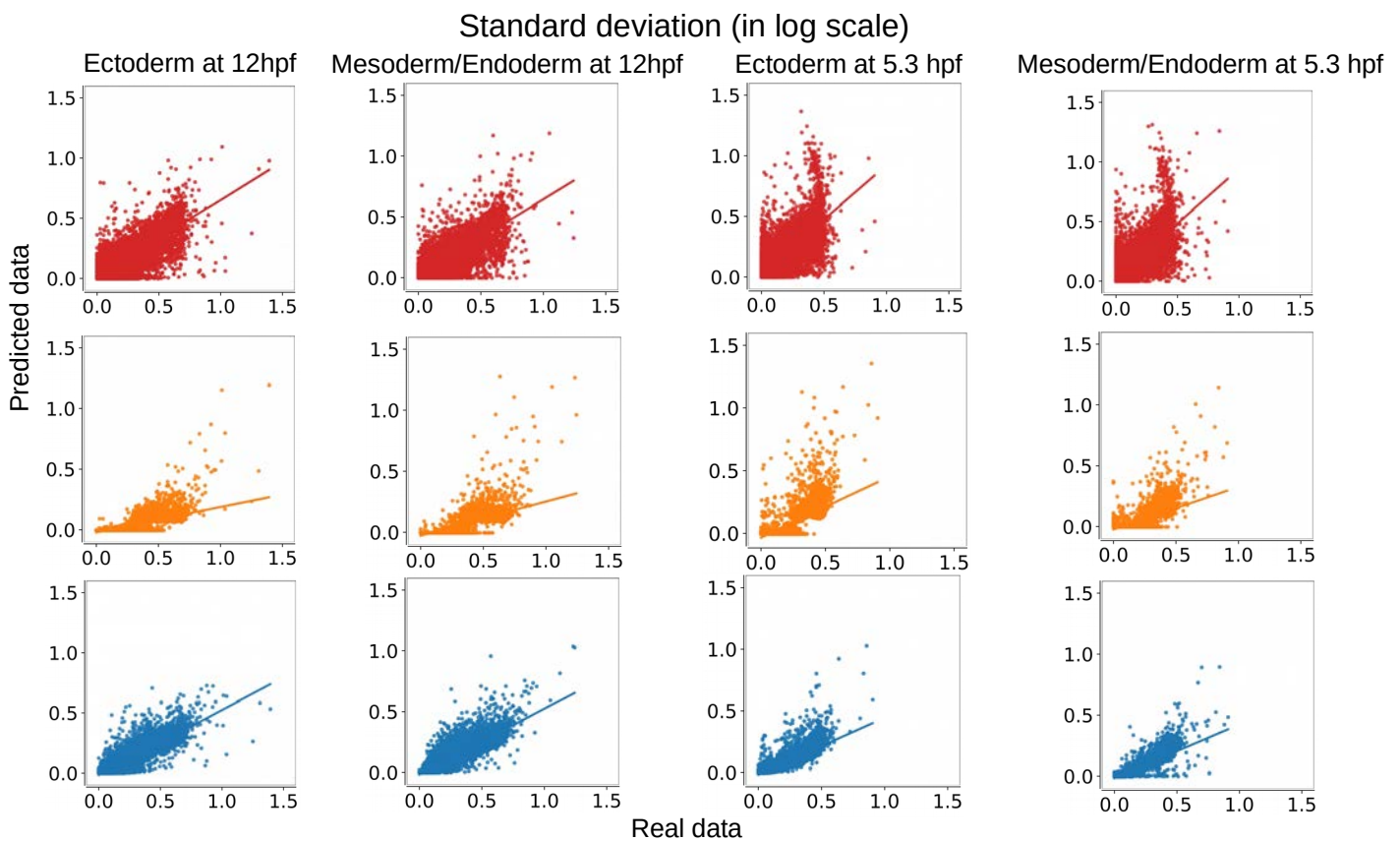
c-d) Zoom-in box plots highlighting variance between repeats.

Data from Farrell et al., 2018⁷, GEO accession GSE106587.

a



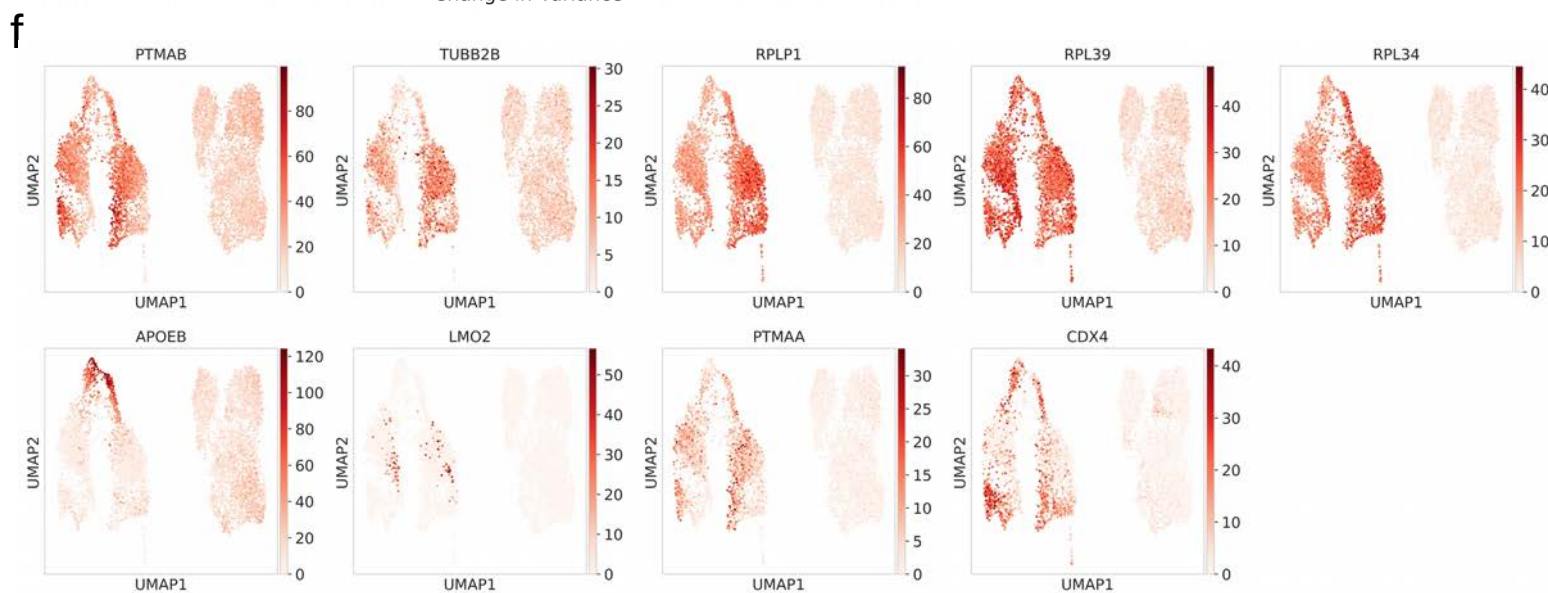
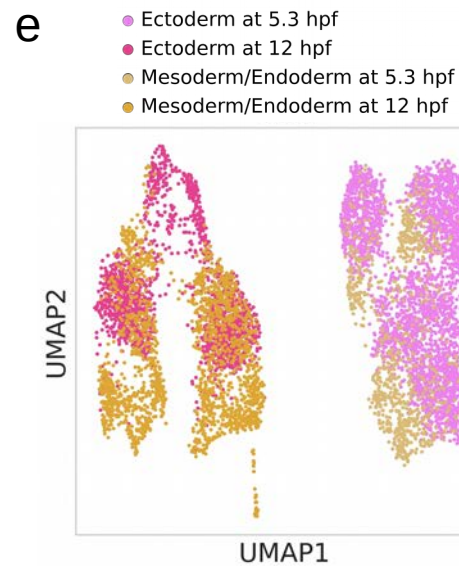
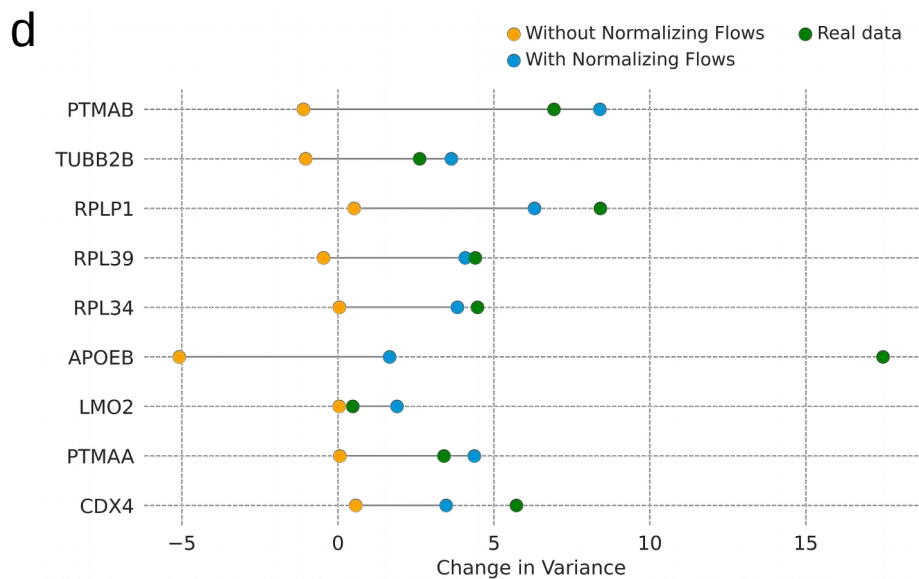
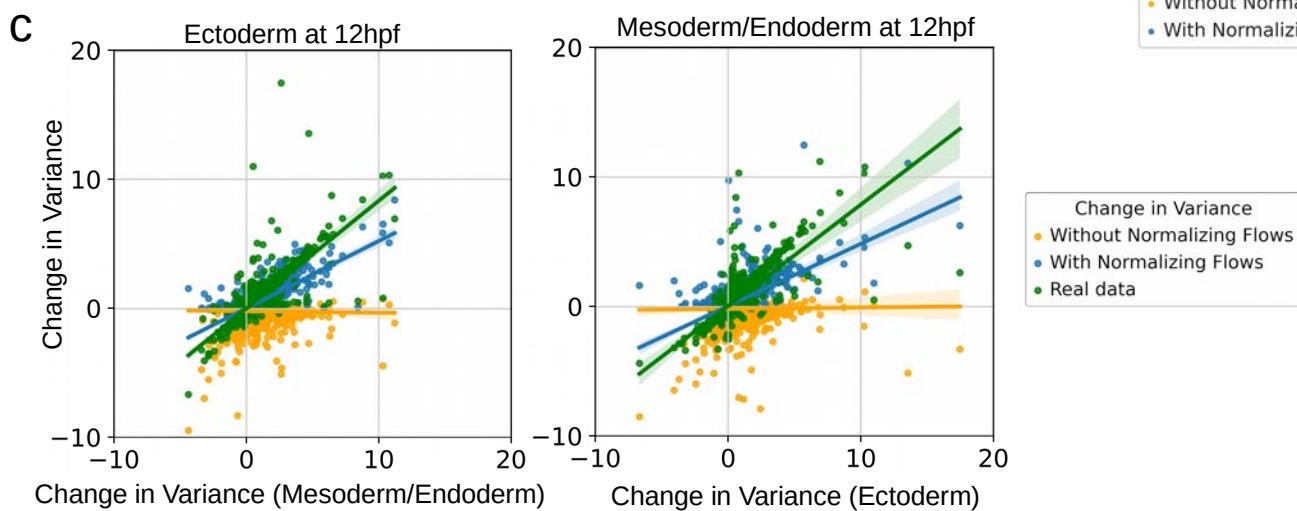
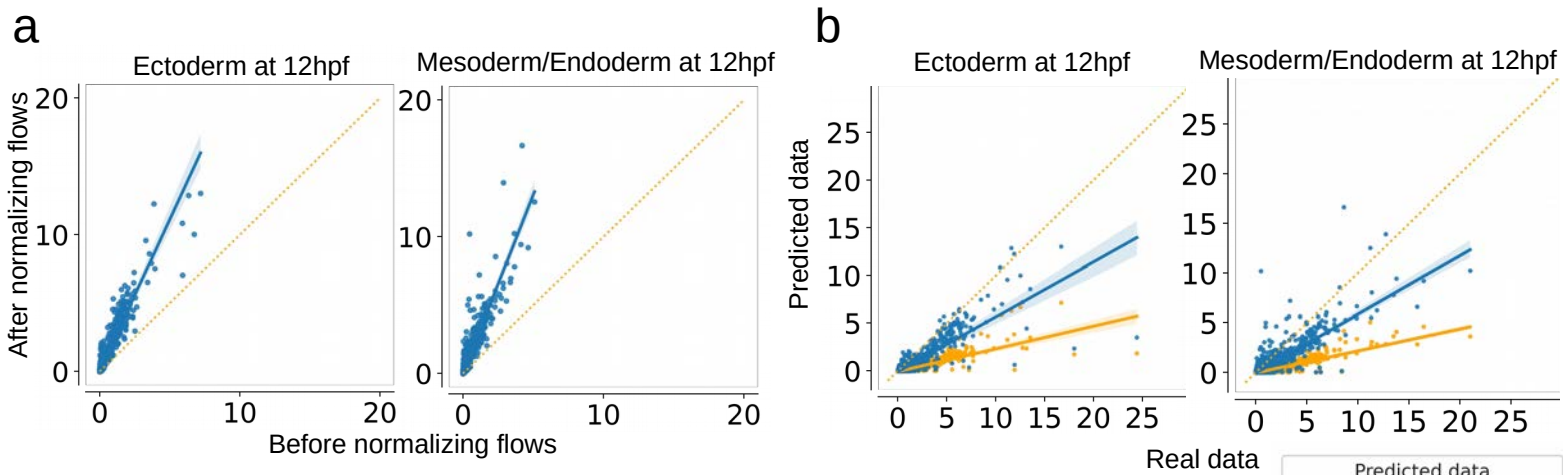
b



Supplementary Figure 5 | Comparison of prediction algorithms (vector arithmetic (VA), variational autoencoders with vector arithmetic (VAE+VA), and mmd-variational autoencoders with flows (mVAE+flows)) in the developing zebrafish (related to Fig. 2).

a-b) Scatter plots of predicted mean gene expression and expression variability across all genes in $\log(x + 1)$ scale with prediction algorithms in the developing zebrafish.

Data from Farrell et al., 2018⁷, GEO accession GSE106587.

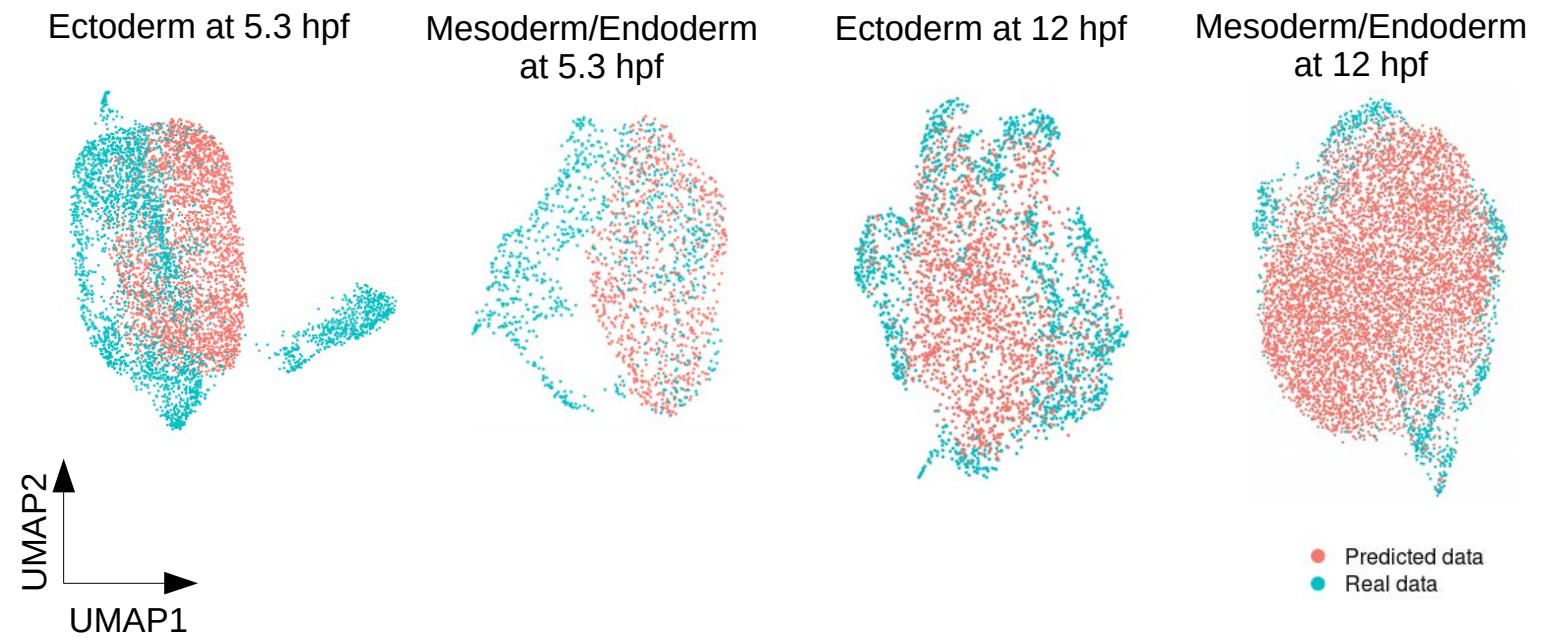


Supplementary Figure 6 | Normalizing flows increase predictive value at high temporal gain of gene variance (related to Fig. 2).

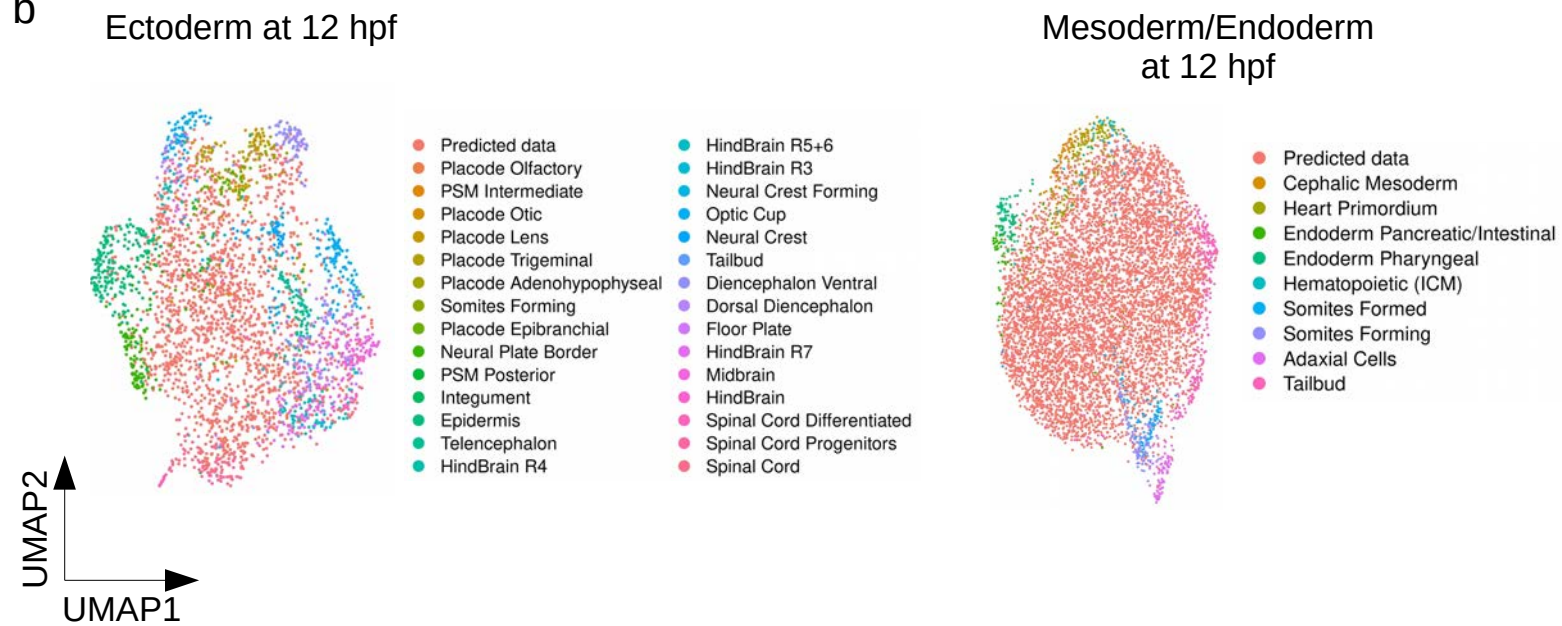
- a)** Comparison of predicted gene expression variance with and without normalizing flows.
- b)** Comparison of predicted variance with and without normalizing flows to real variance.
- c)** Change of expression variance in real and predicted data compared to the training data (in log scale). **a-c)**: predicting to 12hpf ectoderm (left) and 12hp mesoderm/endoderm (right).
- d)** Change of expression variance in real and predicted data compared to the training data for selected genes with high variance gain (in log scale).
- e)** UMAP of test and training data.
- f)** Expression of genes selected in **d)** on the UMAP shown in **e)**.

Data from Farrell et al., 2018⁷, GEO accession GSE106587.

a



b



Supplementary Figure 7 | DCP generates realistic single-cell transcriptomes but does not accurately predict cell type clusters at 12 hpf (related to Fig. 2).

a-b) UMAP of integrated real and predicted data of zebrafish ectoderm and mesendoderm at 5.3 hpf and 12 hpf.

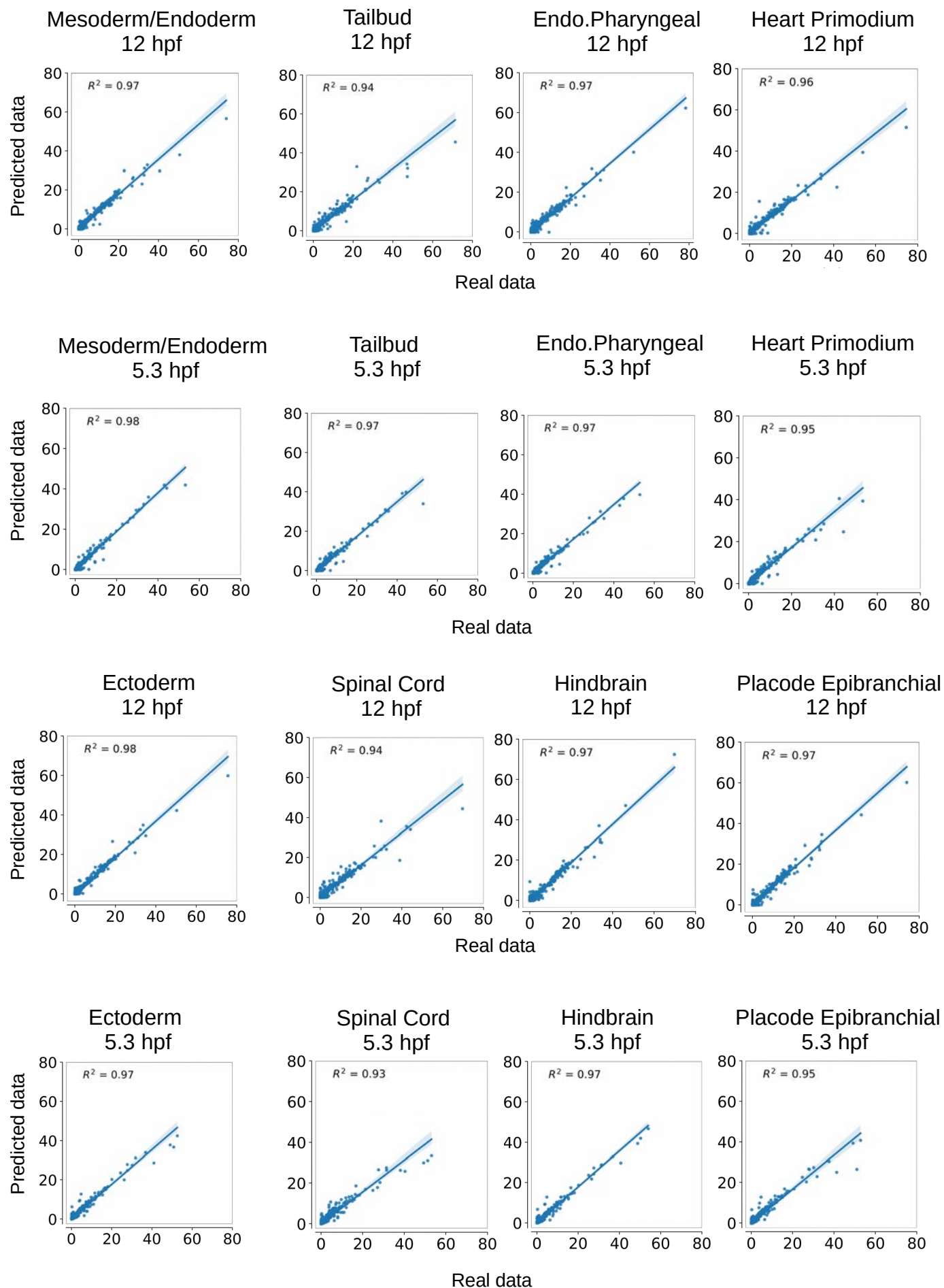
Data from Farrell et al., 2018⁷, GEO accession GSE106587.

a

Mean gene expression

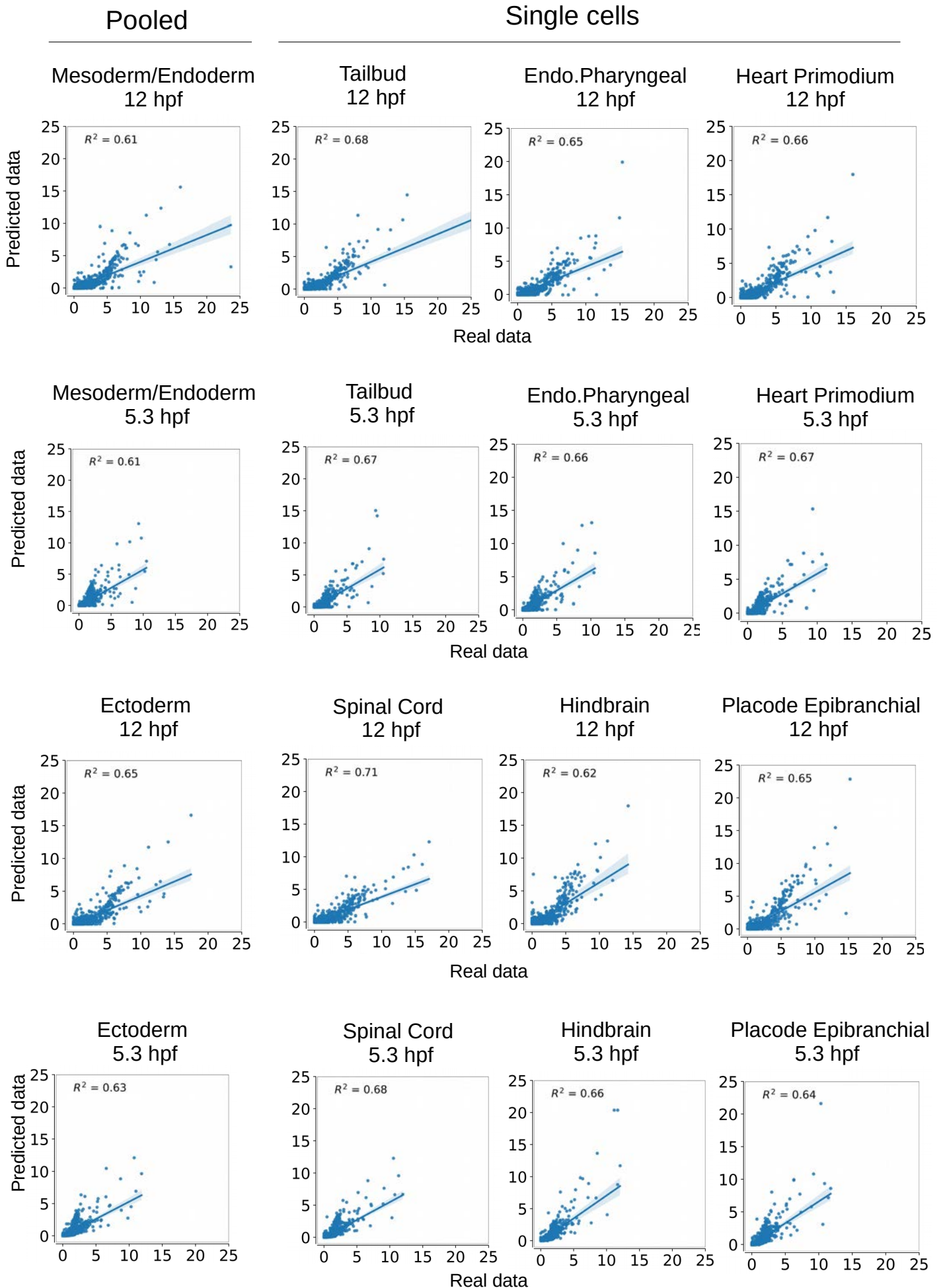
Pooled

Single cells



b

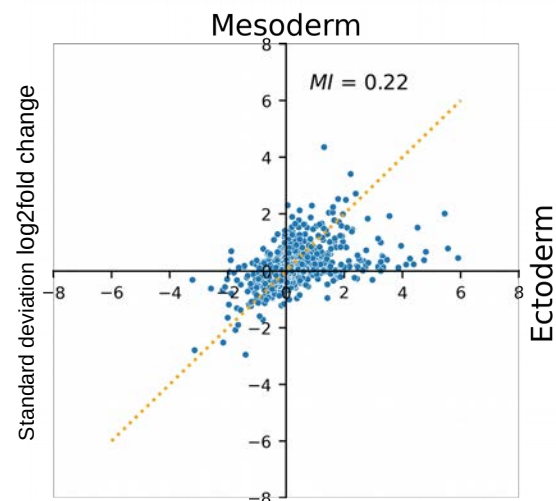
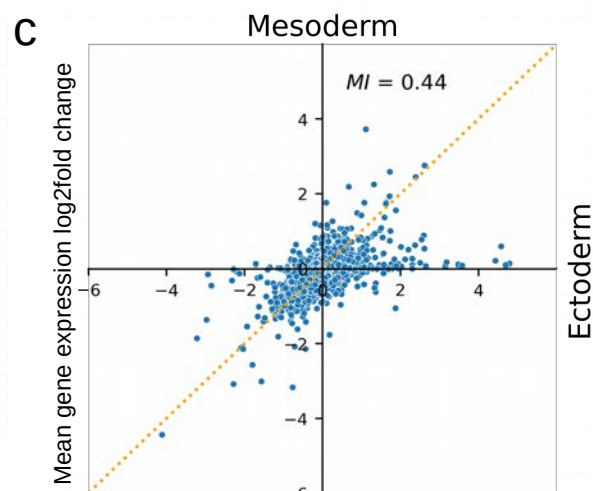
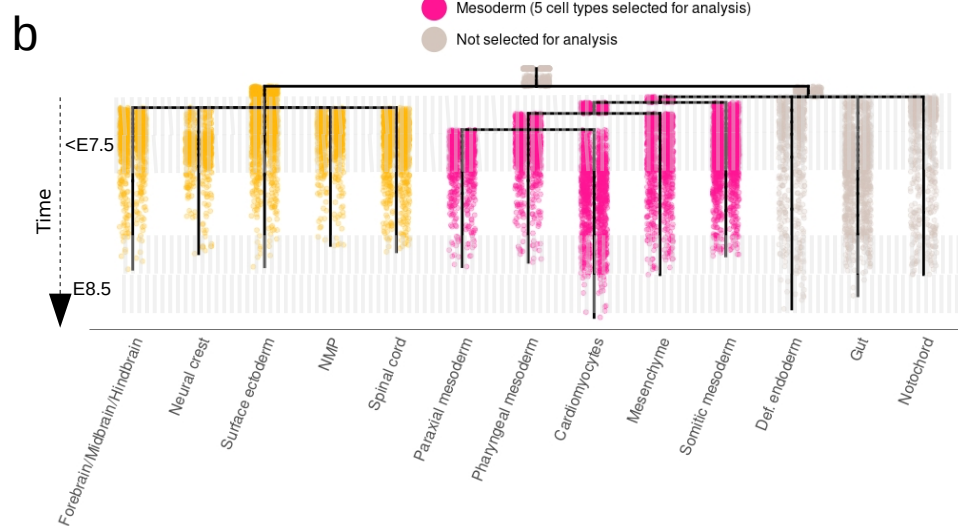
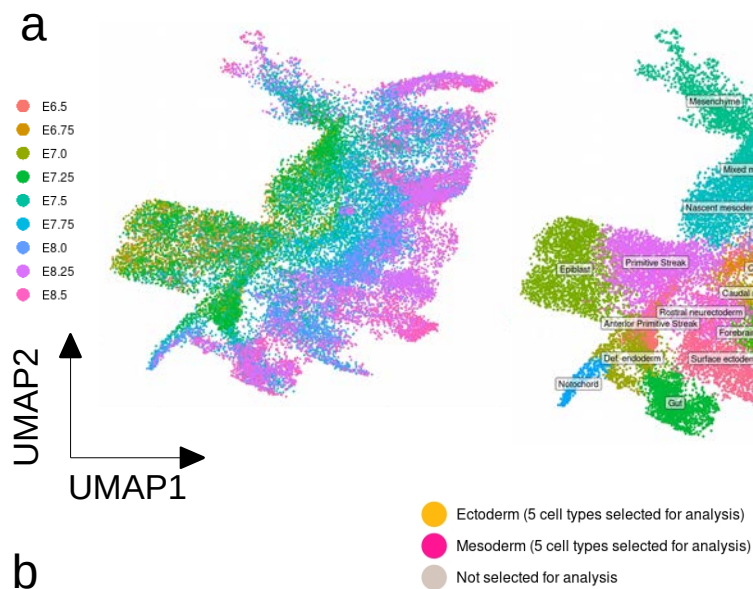
Standard deviation across all genes



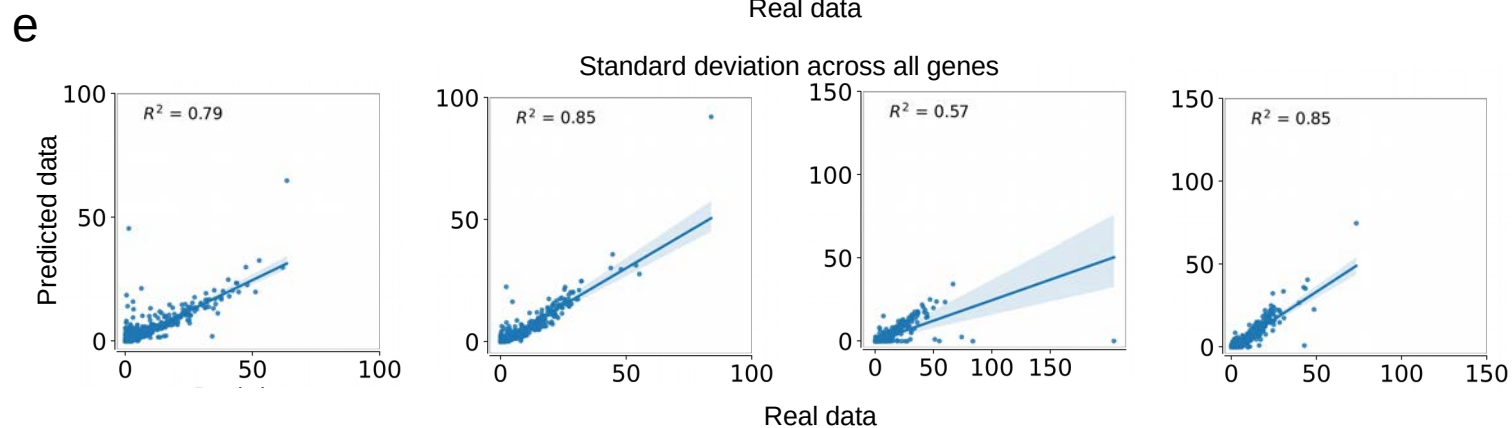
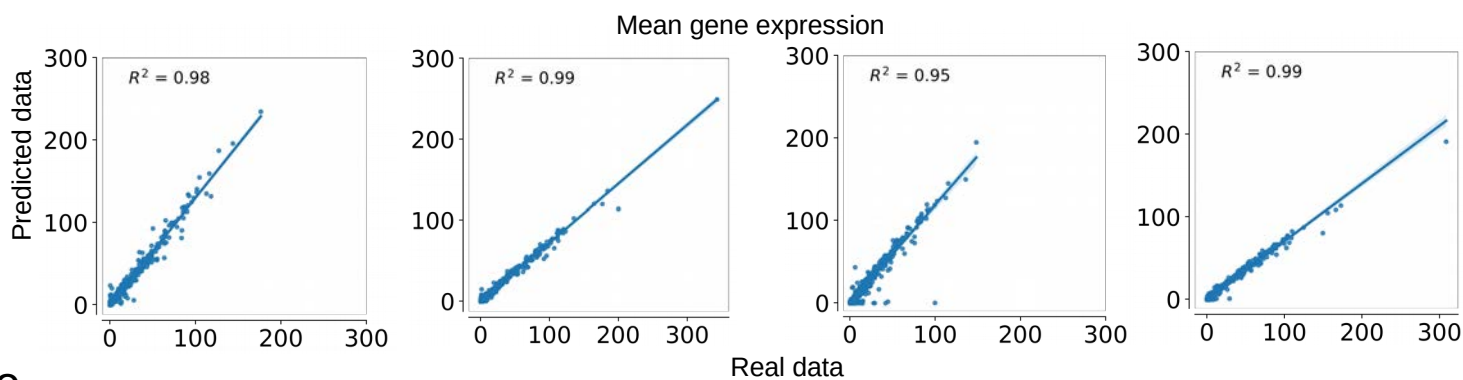
Supplementary Figure 8 | DCP prediction of pooled and single cell types in zebrafish development (related to Fig. 2).

a-b) Mean gene expression and variability were estimated for single cell types and compared to pooled cell types. For this analysis we upsampled the number of cells for three ectodermal cell types (spinal cord, hindbrain, placode epibranchial) and three mesendodermal cell types (tailbud, endo.pharyngeal, heart primordium) to 500 cells. We trained and tested each set of paired ectoderm-mesendoderm cells in both pooled and single type dataset format.

Data from Farrell et al., 2018⁷, GEO accession GSE106587.

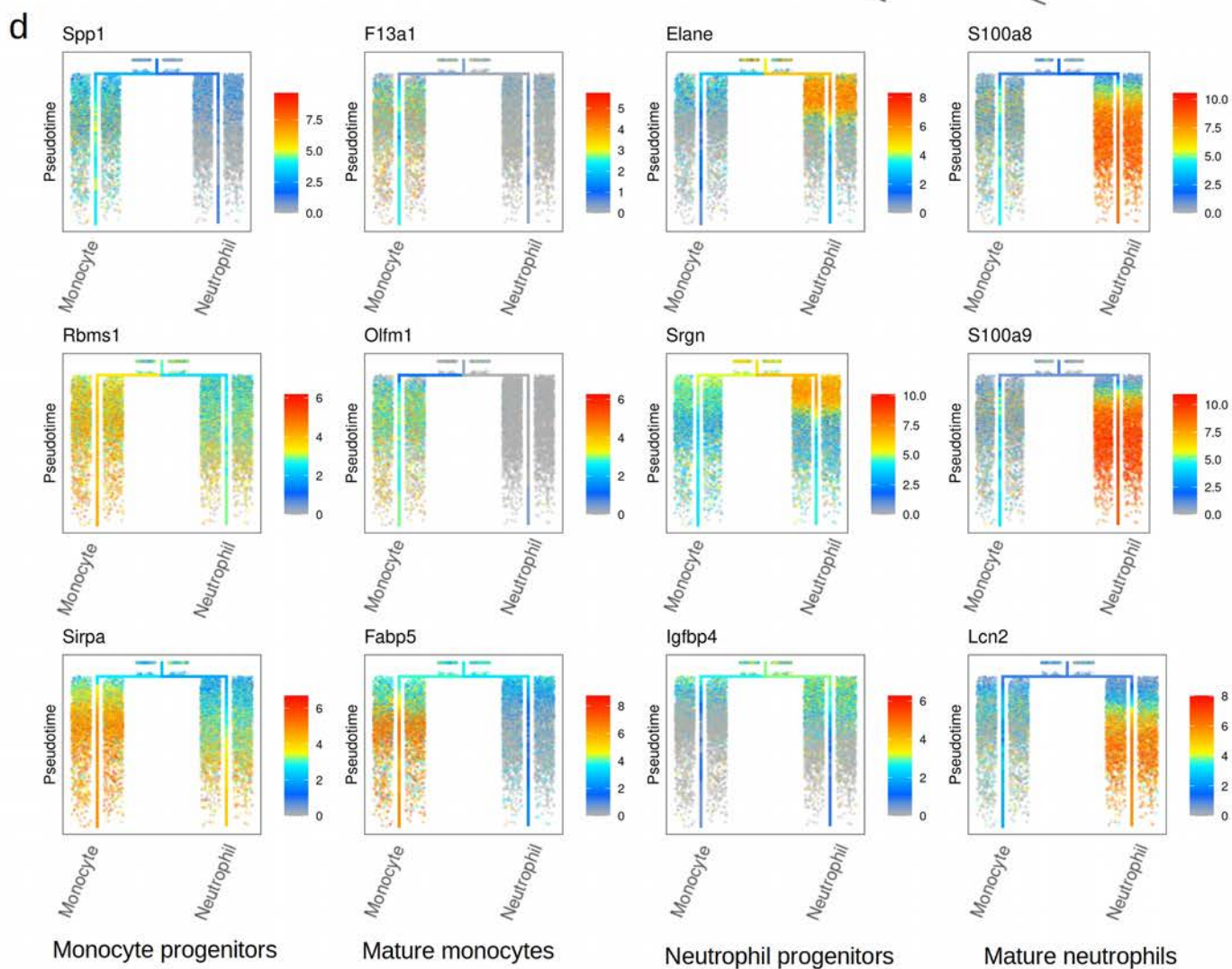
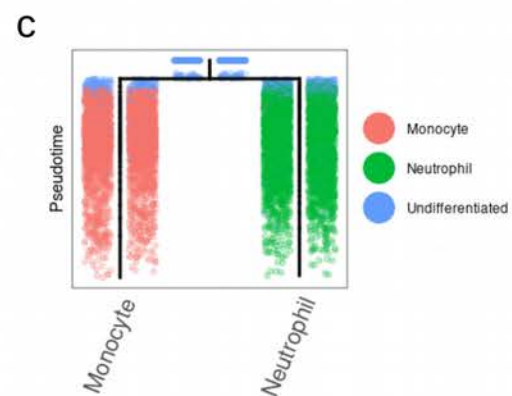
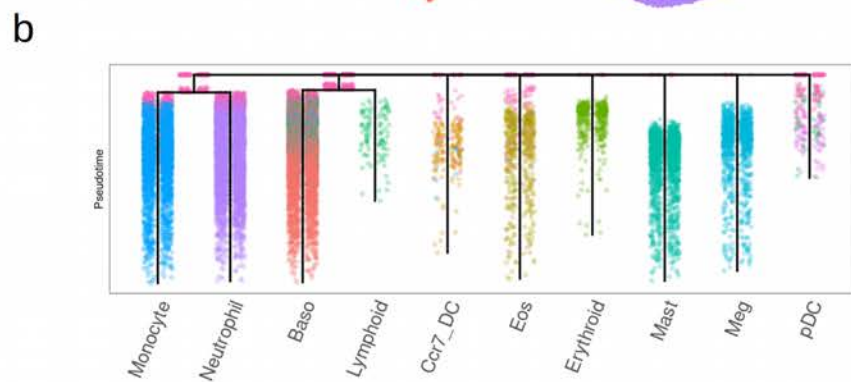
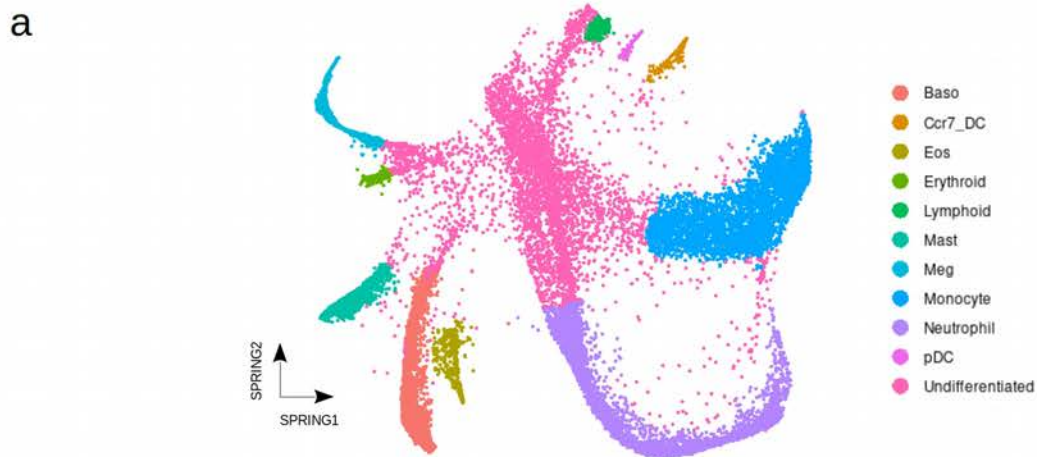


d Forward Ectoderm Backward Ectoderm Forward Mesoderm Backward Mesoderm



Supplementary Figure 9 | Mouse gastrulation dataset (related to Fig. 2).

- a)** UMAP representation of the mouse gastrulation dataset, indicating embryonic day (E). Reproduced from original publication after excluding blood and extraembryonic cells¹⁰.
 - b)** Transcriptome-based tree of same dataset constructed using URD. Similar to our analysis of the zebrafish development dataset, we selected five ectodermal and five mesodermal cell types for the analysis (forward and backward predictions).
 - c)** Mutual information of log2 fold changes between ectoderm and mesoderm for mean and standard deviation of gene expression.
 - d)** Correlation between real and predicted data for mean gene expression.
 - e)** Correlation between real and predicted data for standard deviation of gene expression.
- Data from Pijuan-Sala et al., 2019¹⁰ (ArrayExpress accession E-MTAB-6967).



Supplementary Figure 10 | Mouse hematopoiesis dataset (related to Fig. 3).

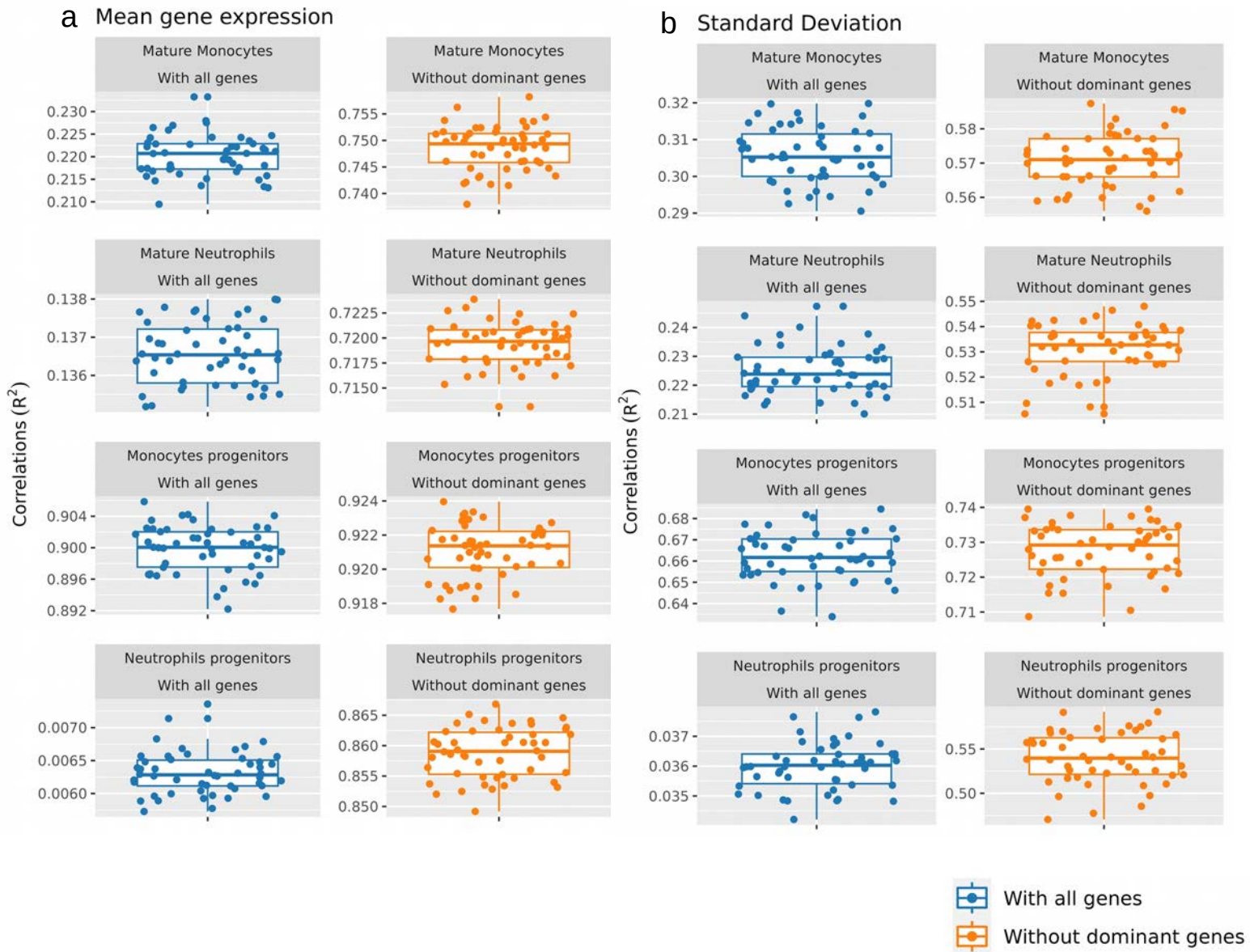
a) t-SNE of dataset.

b) URD of dataset showing branching lineage structure.

c) Subset used in our study: monocytes and neutrophils in a continuum between progenitor and mature states.

d) Representation of selected marker genes to validate that the URD analysis correctly separates mature and progenitor states of monocytes and neutrophils. Color scale represents gene expression in log scale.

Data from Weinreb et al., 2020¹⁸ (GEO accession GSE140802).



Supplementary Figure 11 | Uncertainty in DCP prediction of mature and progenitor cell states of monocyte and neutrophil cell lineages of mouse hematopoiesis (related to Fig. 3).

a-b) Box plots representing distribution of correlation values of predicted mean gene expression and expression variability with all genes and without dominant genes.

Data from Weinreb et al., 2020¹⁸ (GEO accession GSE140802).

Mature Monocytes

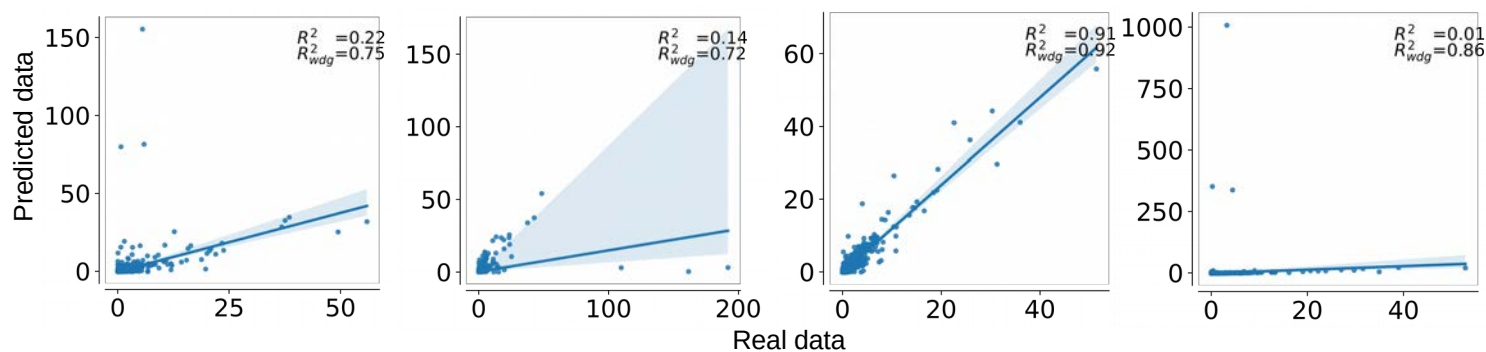
Mature Neutrophils

Monocytes Progenitors

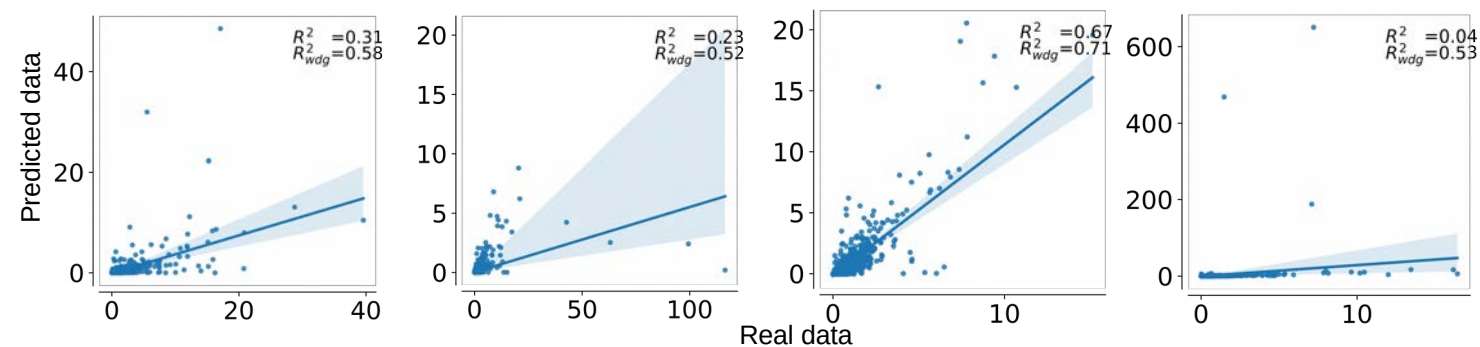
Neutrophils Progenitors

a

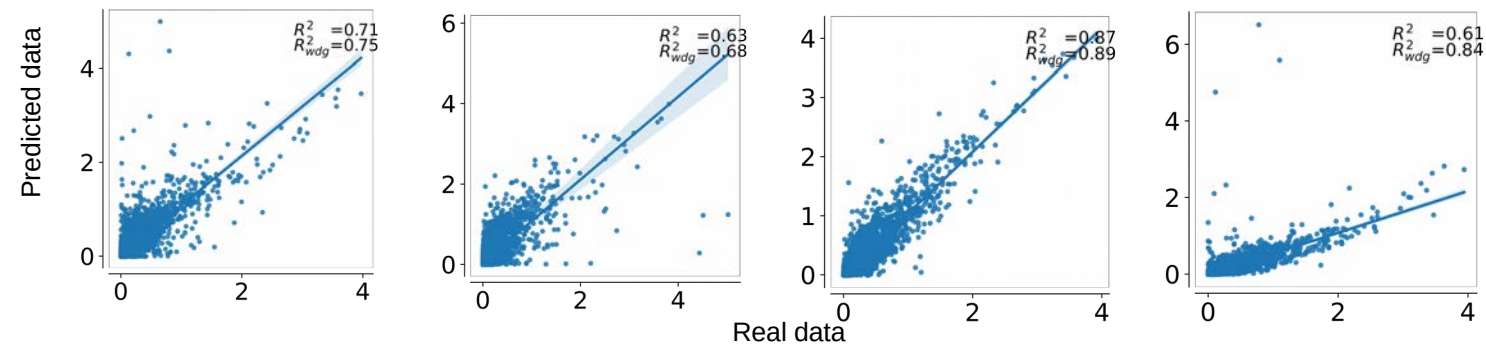
Mean gene expression (in normalized scale)

**b**

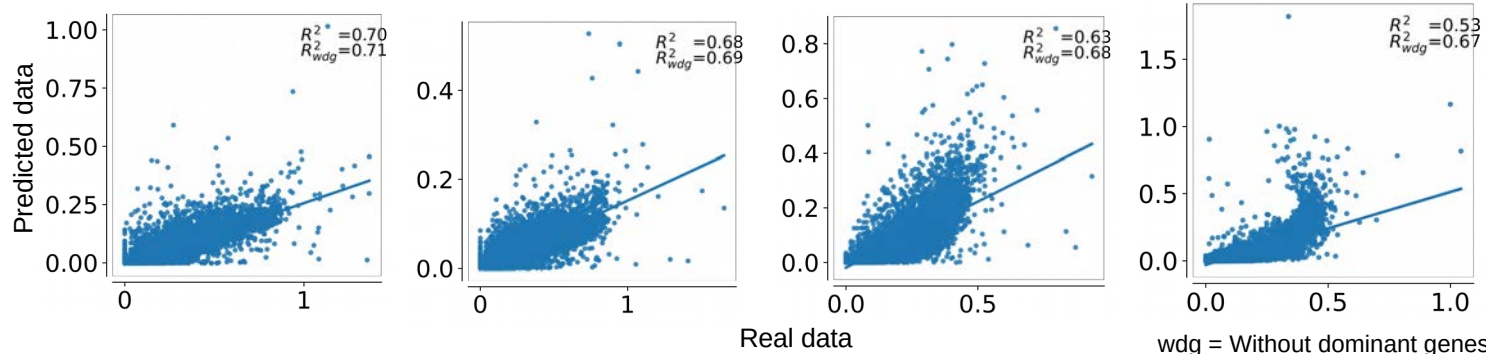
Standard deviation across all genes (in normalized scale)

**c**

Mean gene expression (in log scale)

**d**

Standard deviation across all genes (in log scale)

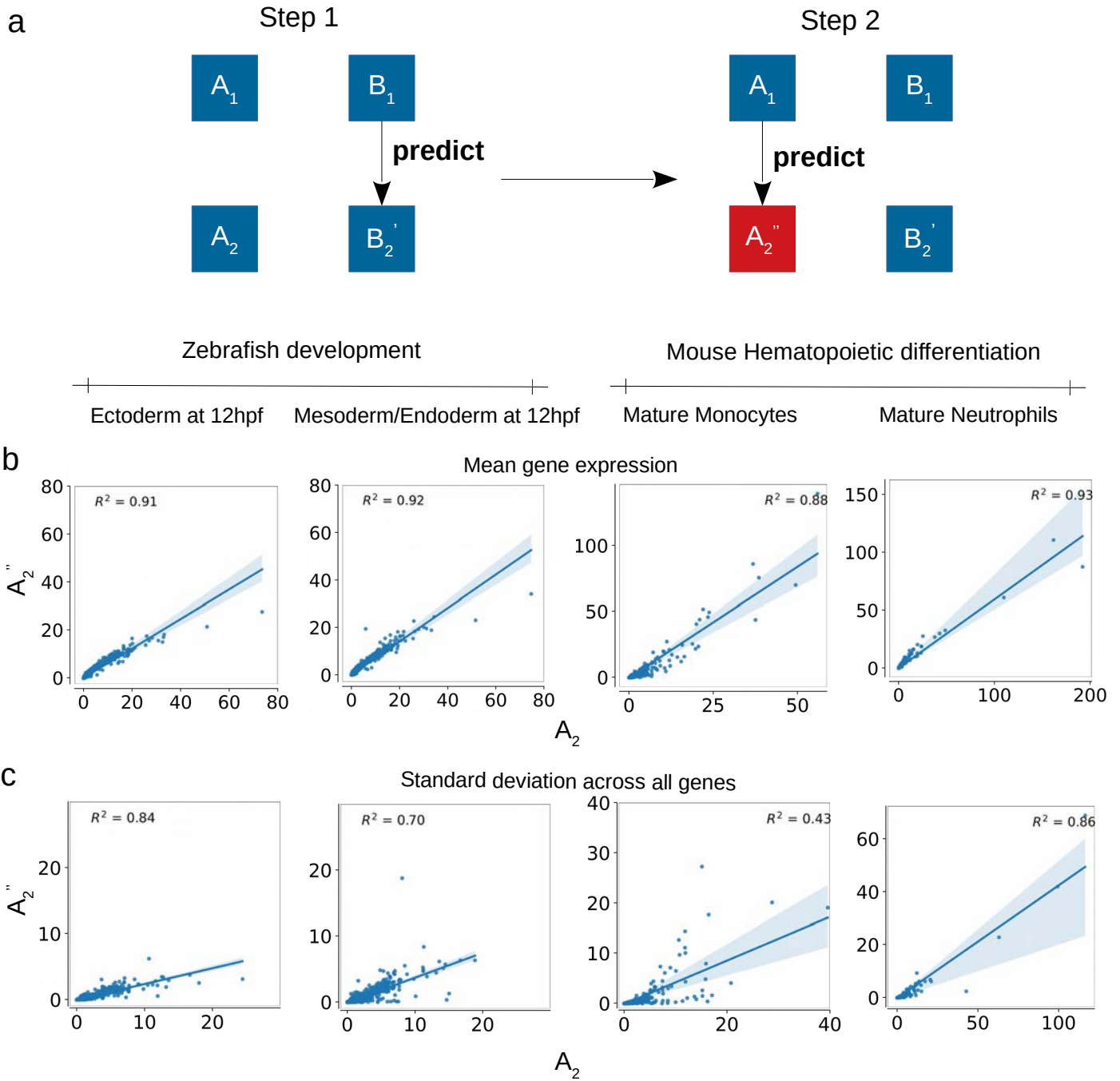


wdg = Without dominant genes

Supplementary Figure 12 | DCP prediction of mature and progenitor cell states of monocyte and neutrophil cell lineages of mouse hematopoiesis (related to Fig. 3).

a-d) Scatter plots of predicted mean gene expression and expression variability across all genes, **a,b)** normalized scale, **c,d)** $\log(x + 1)$ scale.

R^2 and R^2_{wdg} are correlation coefficients calculated with all genes and without dominant genes, respectively. Dominant genes are colored black, other genes are colored blue. Fit lines are lineage regressions with zero intercept. Data from Weinreb et al., 2020¹⁸ (GEO accession GSE140802).

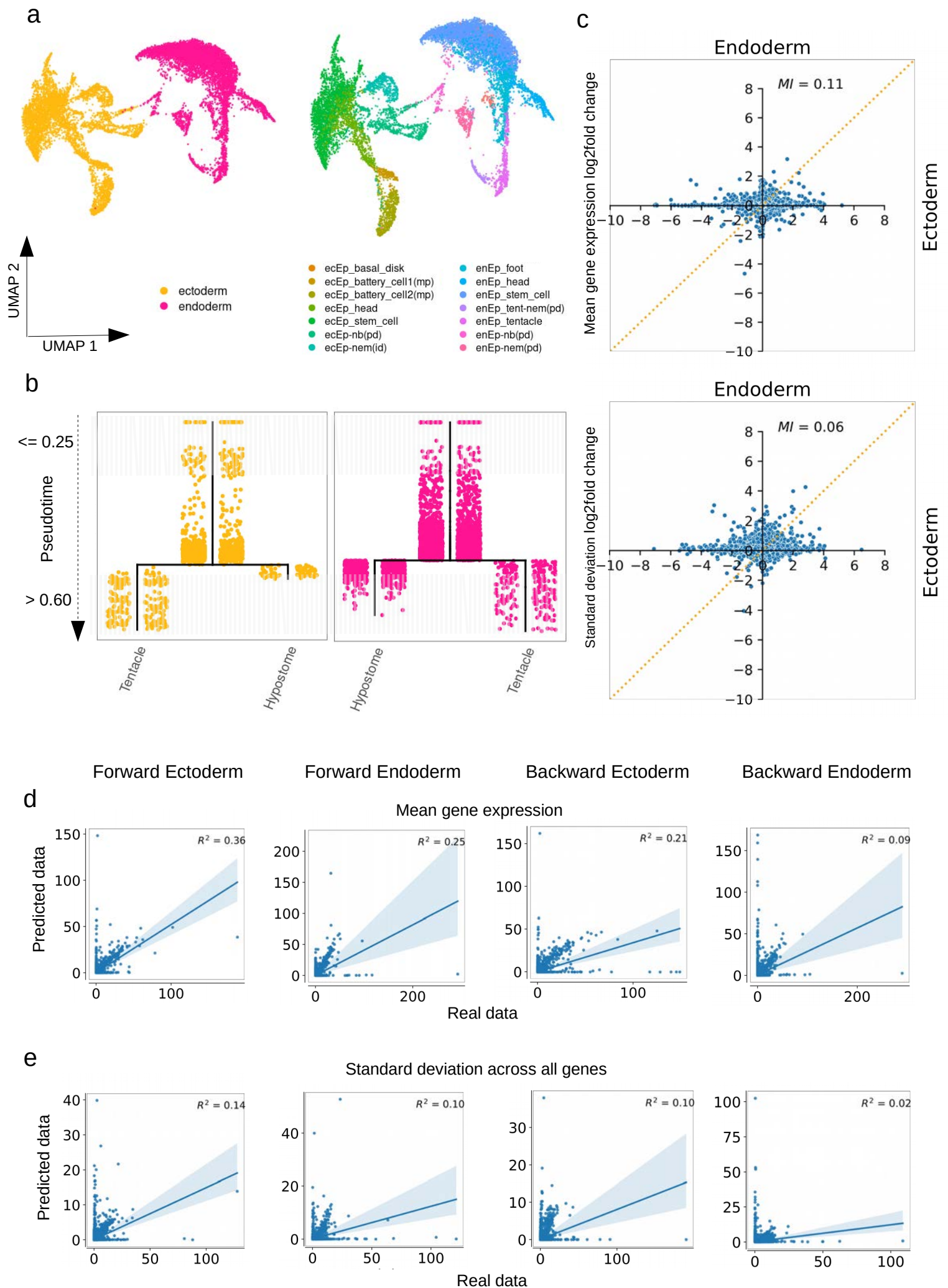


Supplementary Figure 13 | Successive predictions assess the information content of the latent space (related to Fig. 3).

a) To generate a double prediction, we first predict B_2' from A_1 , A_3 and B_1 . We then generate a new prediction for A_2 , A_2'' , from A_1 , B_2 and B_2' .

b-c) Zebrafish (left) and mouse hematopoiesis (right) comparison of double predictions with ground truth on mean gene expression **(b)** and expression standard deviation **(c)**.

Data from Farrell et al., 2018⁷, GEO accession GSE106587 and Weinreb et al., 2020¹⁸ (GEO accession GSE140802).



Supplementary Figure 14 | Hydra stem cell differentiation dataset (related to Fig. 3).

- a)** UMAP representation of ectodermal and endodermal cells from the Hydra dataset³⁰.
 - b)** Transcriptome-based tree of same dataset constructed using URD.
 - c)** Mutual information of log₂ fold changes between ectoderm and endoderm for mean and standard deviation of gene expression.
 - d)** Correlation between real and predicted data for mean gene expression.
 - e)** Correlation between real and predicted data for standard deviation of gene expression.
- Data from Siebert et al., 2019³⁰ (GEO accession GSE121617).