

# Model-based prediction of spatial gene expression via generative linear mapping

Yasushi Okochi<sup>a,b†</sup>, Shunta Sakaguchi<sup>c†</sup>, Ken Nakae<sup>d</sup>, Takefumi Kondo<sup>c,e</sup>, Honda Naoki<sup>a,f\*</sup>

<sup>a</sup> Laboratory for Theoretical Biology, Graduate School of Biostudies, Kyoto University, Sakyo, Kyoto, Kyoto, Japan

<sup>b</sup> Faculty of Medicine, Kyoto University, Sakyo, Kyoto, Kyoto, Japan

<sup>c</sup> Laboratory for Cell Recognition and Pattern Formation, Graduate School of Biostudies, Kyoto University, Sakyo, Kyoto, Kyoto, Japan

<sup>d</sup> Graduate School of Informatics, Kyoto University, Sakyo, Kyoto, Kyoto, Japan

<sup>e</sup> The Keihanshin Consortium for Fostering the Next Generation of Global Leaders in Research (K-CONNEX), Sakyo, Kyoto, Kyoto, Japan

<sup>f</sup> Research Center for Dynamic Living Systems, Kyoto University, Sakyo, Kyoto, Kyoto, Japan

† These authors contributed equally to this manuscript.

## \* Corresponding author.

Graduate School of Biostudies, Kyoto University, Yoshidakonoecho, Sakyo-ku, Kyoto, Kyoto 606-8315, Japan.

Tel.: +81-75-753-9450

E-mail: honda.naoki.4v@kyoto-u.ac.jp

## Abstract

Decoding spatial transcriptomes from single-cell RNA sequencing (scRNA-seq) data has become a fundamental technique for understanding multicellular systems; however, existing computational methods lack both accuracy and biological interpretability due to their model-free frameworks. Here, we introduced Perler, a model-based method to integrate scRNA-seq data with reference *in situ* hybridization (ISH) data. To calibrate differences between these datasets, we developed a biologically interpretable model that uses generative linear mapping based on a Gaussian-mixture model using the Expectation-Maximization algorithm. Perler accurately predicted the spatial gene expression of *Drosophila* embryos, zebrafish embryos, mammalian liver, and mouse visual cortex from scRNA-seq data. Furthermore, the reconstructed transcriptomes did not over-fit the ISH data and preserved the timing information of the scRNA-seq data. These results demonstrated the generalizability of Perler for dataset integration, thereby providing a biologically interpretable framework for accurate reconstruction of spatial transcriptomes in any multicellular system.

## 1 Introduction

2 Genes are heterogeneously expressed in multicellular systems, and their spatial profiles are tightly linked to  
3 biological functions. In developing embryos, spatial gene-expression patterns are responsible for coordinated  
4 cell behavior (e.g., differentiation and deformation) that regulates morphogenesis<sup>1</sup>. Additionally, within organ  
5 tissues, cells at different locations play different roles in organ function based on their gene-expression  
6 patterns<sup>2</sup>. Thus, identification of spatial genome-wide gene-expression profiles is key to understanding the  
7 functions of various multicellular systems. *In situ* hybridization (ISH) has been widely used to visualize spatial  
8 profiles of gene expression; however, application of this method is generally limited to only small numbers of  
9 genes. By contrast, the single-cell RNA sequencing (scRNA-seq) method developed during the previous  
10 decade has enabled measurement of genome-wide gene-expression profiles in tissues at the single-cell level<sup>3</sup>.  
11 However, this method requires tissue dissociation, which leads to loss of spatial information for the original  
12 cells.

13 To compensate for the lost spatial information, new computational approaches have emerged (Seurat  
14 v.1<sup>4</sup>, DistMap<sup>5</sup>, Achim et al.<sup>6</sup>, Halpern et al.<sup>7</sup>), enabling reconstruction of genome-wide spatial expression  
15 profiles from scRNA-seq data by integrating existing ISH data as a spatial reference map *in silico*. However,  
16 their methods require binarization of gene-expression data<sup>8</sup>, which leads to unsatisfactory accuracy, or tissue-  
17 specific modelling, which leads difficulty in application to other systems. Recently, the seminal methods Seurat  
18 (v.3)<sup>9</sup> and Liger<sup>10</sup> were developed to address gene-expression data as continuous variables in a non-tissue-  
19 specific manner. These methods match the distributions of ISH and scRNA-seq data points by using  
20 dimensionality reductions [e.g., canonical correlation analysis (CCA)]<sup>11</sup> and integrative non-negative matrix  
21 factorization (iNMF)<sup>12</sup>, followed by mapping the scRNA-seq data points to the nearest ISH data points  
22 according to Euclidean distance using Nearest-Neighbor (NN) methods (e.g., k-NN<sup>13</sup> and mutual NN<sup>14</sup>).  
23 However, a major issue is that the flexibility of the methods allow mapping of ISH data to scRNA-seq data  
24 without any models of the underlying scRNA-seq data structure. Specifically, these methods do not account  
25 for difference in gene-expression noise associated with each gene. Given this model-free property, these  
26 methods are dependent upon nonlinear NN mapping, which innately causes over-fitting to the reference ISH  
27 data.

28 To address these issues, we propose a novel model-based computational method for probabilistic  
29 embryo reconstruction by linear evaluation of scRNASEq (Perler), which reconstructs spatial gene-expression  
30 profiles via generative linear modeling in a biologically interpretable framework. Perler addresses gene-  
31 expression profiles as continuous variables and models generative linear mapping from ISH data points into  
32 the scRNA-seq space. To estimate parameters of the linear mapping, we developed a method based on the  
33 Expectation–Maximization (EM) algorithm<sup>13</sup>. Using the estimated parameters, we also propose an  
34 optimization method to infer spatial information of scRNA-seq data within a tissue sample. We applied this  
35 method to existing *Drosophila* scRNA-seq data<sup>5</sup> and successfully reconstructed spatial gene-expression  
36 profiles in *Drosophila* early embryos that were more accurate than those generated using another spatial  
37 reconstruction method (DistMap<sup>5</sup>). Additionally, we showed that Perler can reconstruct a spatial gene-

1 expression pattern that could not be fully predicted using previous methods, including Seurat (v.3), Liger, and  
2 DistMap. Further analysis revealed that Perler was able to preserve the timing information of the scRNA-seq  
3 data without over-fitting to the reference ISH data. Furthermore, we demonstrated that this method accurately  
4 predicted spatial gene-expression profiles in early zebrafish embryos<sup>4</sup>, the mammalian liver<sup>7</sup>, and the mouse  
5 visual cortex<sup>15,16</sup>. These findings demonstrate Perler as a robust, generalized framework for predicting spatial  
6 transcriptomes from any type of ISH data for any multicellular system without overfitting to the reference.

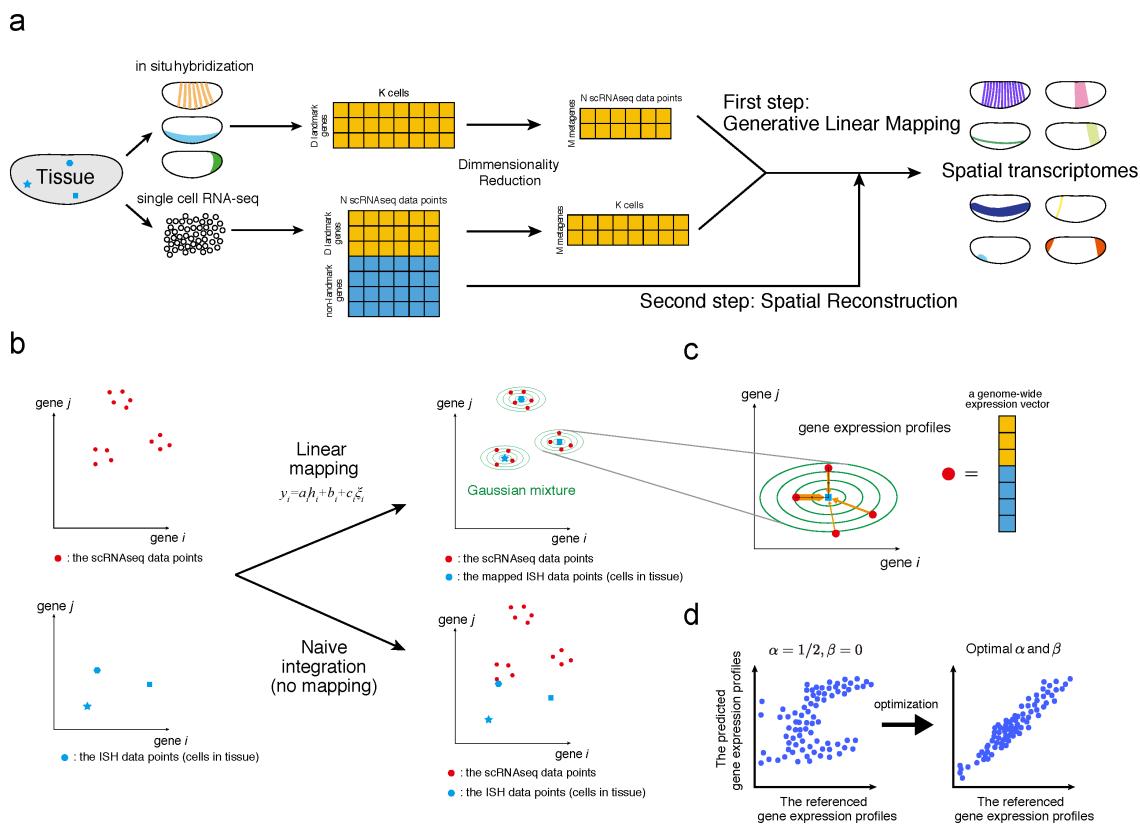
7

8

# 1 Results

## 2 Framework of spatial reconstruction in Perler

3 Perler is a novel computational method for model-based prediction of spatial genome-wide expression profiles  
 4 from scRNA-seq data that works by referencing spatial gene-expression profiles measured by ISH (Fig. 1a).  
 5 In general, scRNA-seq data have higher dimensionality (on the order of ~10,000 genes) but does not contain  
 6 information of spatial coordinates in tissues. By contrast, reference ISH data contain expression information  
 7 for  $D$  genes in each cell or tissue subregion, with these referred to as landmark genes (e.g.,  $D = 84$  in  
 8 *Drosophila melanogaster* early embryos) and tagged with spatial coordinates in tissues.



9

10 **Figure 1: Schematic illustration of Perler**

11 (a) Flow of data processing. (b) Generative linear mapping from ISH data to the scRNA-seq space. The left and right  
 12 panels indicate scatter plots in high-dimensional ISH and scRNA-seq spaces. Because ISH data points did not match  
 13 scRNaseq data points (Naive integration) in the absence of mapping, ISH data points were mapped in order to fit the  
 14 scRNA-seq data points the best using the EM algorithm. Blue points indicate ISH data points, red points indicate scRNA-  
 15 seq data points, and green lines indicate contours of the estimated multivariate Gaussian distribution (see Methods). (c)  
 16 Reconstruction/prediction of gene expression by Mahalanobis' metric-based weighting (see Methods). Orange arrows  
 17 indicate weights between scRNA-seq data points to cell  $k$ , and their widths reflect the Mahalanobis' metric-based weights.  
 18 Note that although two scRNaseq data points are equally distant from the ISH data point in the Euclid metric (black line),  
 19 the weights are different. (d) Weight determination. The hyperparameters of the weighting function,  $\alpha$  and  $\beta$ , are  
 20 determined by cross-validation to ensure that the referenced gene-expression profiles correlate well with the predicted  
 21 gene-expression profiles (see Methods). Dots correspond to cells in tissue. The left and right panels indicate the conceptual  
 22 scatter plots of the expression levels of the genes before ( $\alpha = 1/2, \beta = 0$ ) and after parameter optimization, respectively.

23

1        The Perler procedure involves two steps. The first step estimates a generative linear model-based  
2 mapping function that transforms ISH data into the scRNA-seq space, thereby enabling calculation of pairwise  
3 distances between ISH data and scRNA-seq data (**Fig. 1b**). The second step reconstructs spatial gene-  
4 expression profiles according to the weighted mean of scRNA-seq data, which is optimized by the mapping  
5 function estimated in the first step (**Fig. 1c**).

6        The first step considers gene-specific differences between scRNA-seq and ISH measurements. For  
7 example, we assume that some genes are more or less sensitive to ISH or scRNA-seq and subject to high or  
8 low background signals in the associated data. We account for gene-specific noise intensity, because gene  
9 expression fluctuates over time in a gene-specific manner<sup>17</sup>. These differences in sensitivity, background  
10 signals, and noise intensity can be expressed by linear mapping:

$$y_i = a_i h_i + b_i + c_i \xi_i$$

11        where  $y_i$  and  $h_i$  denote the expression levels of landmark gene  $i$  measured by scRNA-seq and ISH, respectively;  
12  $a_i$ ,  $b_i$ , and  $c_i$  are constant parameters of gene  $i$  and interpreted as the sensitivity coefficient, background  
13 signal, and noise intensity, respectively; and  $\xi_i$  indicates standard Gaussian noise. Note that  $a_i$ ,  $b_i$ , and  $c_i$   
14 are different for each gene, and that these parameter values are unknown. To estimate this linear mapping from  
15 the data, we developed a generative model in which scRNaseq data points are generated/derived from each  
16 cell in the tissue, whose expression is measured by ISH (see Methods). We then derived a parameter-estimation  
17 procedure based on the EM algorithm (see Methods). Using the estimated parameters, a gene-expression  
18 vector for each cell in a given tissue sample measured by ISH can be mapped to the scRNA-seq space, thereby  
19 allowing evaluation of pairwise distances between ISH and scRNA-seq data.

20        The second step reconstructs the spatial gene-expression profile in tissue from scRNA-seq data. We  
21 estimated gene expression of each cell in a tissue sample according to the weighted mean of all scRNA-seq  
22 data points, where the weights were determined by the pairwise distances between cells in tissue samples  
23 measured using ISH and scRNA-seq data points (**Fig. 1c**). For the best prediction, we optimized the  
24 hyperparameters of the weighting function to ensure that the predicted and referenced landmark gene-  
25 expression profiles were well-correlated by cross-validation (**Fig. 1d**; see Methods). We then predicted non-  
26 landmark gene expression using the optimized weighting function.

27        Data were preprocessed before the first step. Some landmark genes redundantly exhibit similar  
28 spatial expression patterns, which can lead to biased parameter estimation and cause a loss of mapping ability.  
29 To reduce redundancy in scRNA-seq and ISH data, we performed dimensionality reduction using partial least  
30 squares correlation analysis (PLSC)<sup>18</sup> (see Methods). Each factor in the reduced dimension can be interpreted  
31 as a “metagene”, which is representative among a highly correlated gene cluster, with its coordinate  
32 corresponding to the expression level of the metagene. In Perler, we regarded the metagene-expression level  $i$   
33 (i.e., factor  $i$ ) in the scRNA-seq and ISH spaces as  $y_i$  and  $h_i$  in the equation above (see Methods).

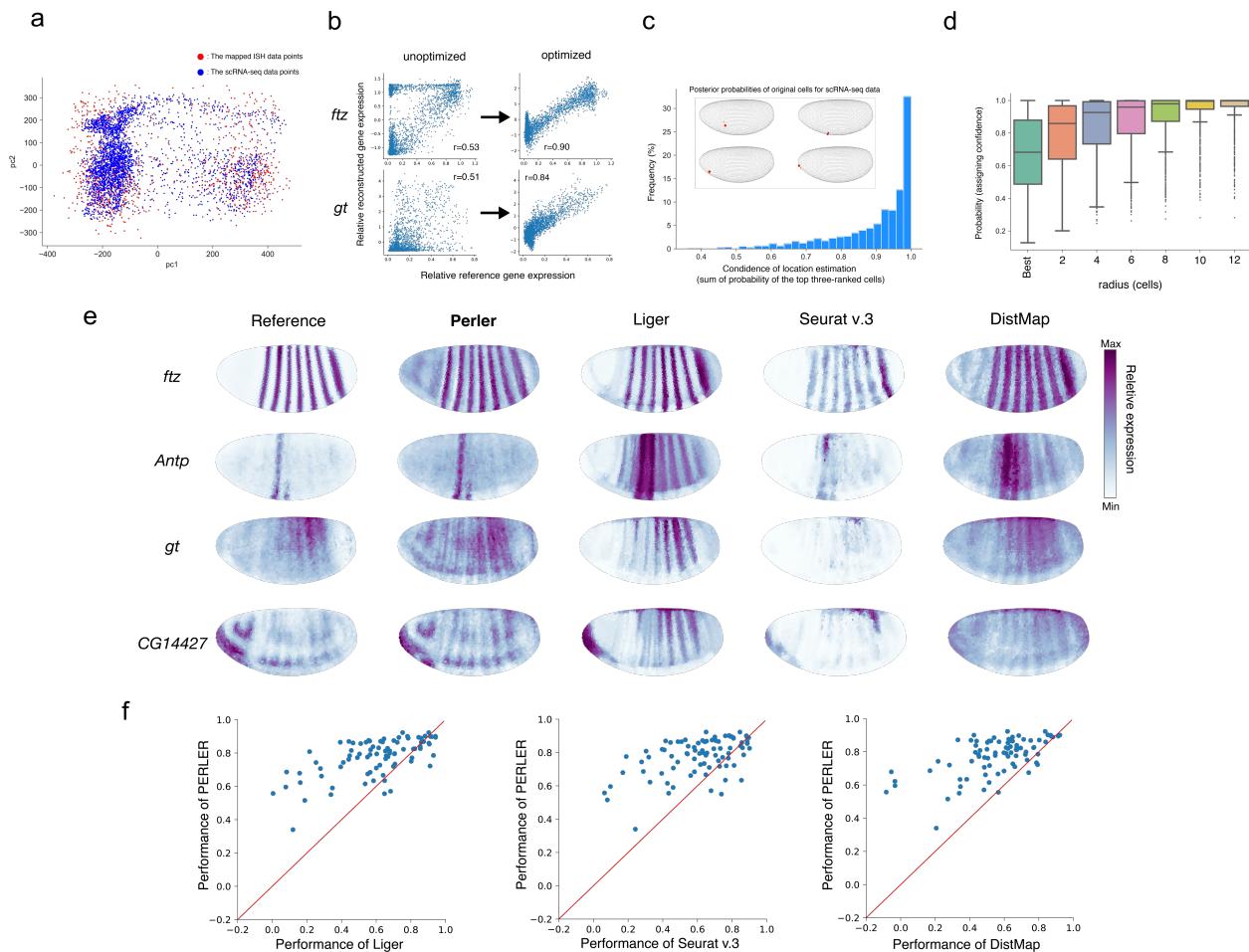
35

### 36 **Model-based mapping between scRNA-seq and ISH data**

37        Previously, Karaïkos *et al.*<sup>5</sup> measured gene expression in individual cells dissociated from early *D.*  
38 *melanogaster* embryos at developmental stage 6 by scRNA-seq, followed by development of a computational

1 method (DistMap) to reconstruct the spatial gene-expression profile of the embryos from the scRNA-seq data.  
2 They used as reference data a spatial gene-expression atlas provided by the Berkeley *Drosophila* Transcription  
3 Network Project (BDTNP)<sup>19,20</sup>, in which the expression of 84 landmark genes was quantitatively measured by  
4 fluorescent (FISH) at single-cell resolution at developmental stage 5.

5



6

## 7 **Figure 2: Generation of spatial gene expression profiles**

8 (a) Scatter plot of mapped gene expression and scRNA-seq observations (Fig. 1b, upper right panel). Principal component  
9 analysis<sup>13</sup> was used to visualize high-dimensional gene expression data into two dimensions. (b) Improved correlation  
10 between predicted and referenced data in the scRNA-seq space by optimizing the weighting function. (c) Histogram of the  
11 assigned confidence calculated as the posterior probabilities of the top three cells for each scRNA-seq data point. Perler  
12 assigned the majority of scRNA-seq data points (82.7%) to the top three cells with high confidence (>0.8). Examples of  
13 estimated original positions of each scRNA-seq data point (inset). Embryos are colored according to the posterior  
14 probabilities for the scRNA-seq data points. Point sizes indicate the magnitude of each posterior probability. Points with  
15 posterior probabilities below 0.001 were omitted. (d) Boxplot of the assigned specificity calculated as the posterior  
16 probabilities of circular regions for each scRNA-seq data point according to radius, with the center of each region  
17 representing the best assigned location for each data point. The boxplot has whiskers with a maximum 1.5 interquartile  
18 range, with black points indicating outliers. The radius was calculated by path length on the k-NN graph comprising all cells  
19 in the tissue ( $k = 6$ ). (e) Reconstructions of the landmark genes by Perler, Liger, Seurat (v.3), and DistMap. (f) Comparison  
20 of Perler reconstruction performance with Liger (left), Seurat (v.3) (middle), and DistMap (right). Each dot indicates the  
21 reconstruction accuracies for each gene by Perler and other methods.

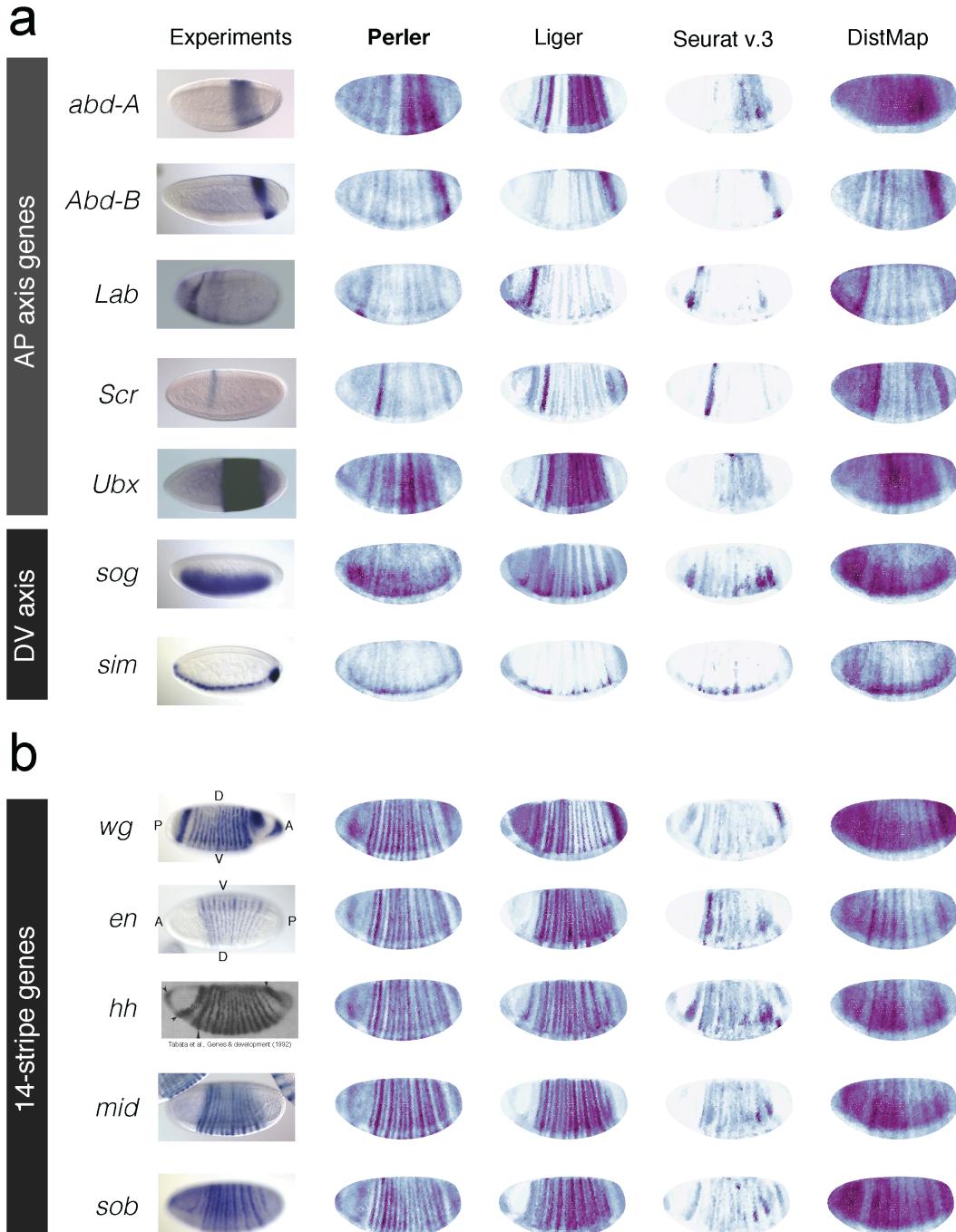
22

1 In the present study, we applied Perler to the same scRNA-seq dataset and used the 84 landmark  
2 genes from the BDTNP atlas as the spatial reference map. We then predicted the spatial gene-expression  
3 profiles for 8840 non-landmark genes. To compare Perler results with those of DistMap, we used the same  
4 normalization methods for the scRNA-seq dataset as the previous study<sup>5</sup>. For preprocessing, we manually  
5 extracted 60 metagenes as non-redundant clusters of the landmark genes by dimensionality reduction and  
6 estimated the parameters of the linear mapping by integrating the scRNA-seq data with the ISH data  
7 (**Supplementary Fig. 1a–c**). The mapped ISH data points according to the linear mapping were distributed  
8 consistently with the scRNA-seq data points (**Fig. 2a**, **Supplementary Fig. 1d**). Additionally, the  
9 reconstructed and referenced gene-expression profiles were well-correlated following optimization of the  
10 hyperparameters (**Fig. 2b**). We then derived a posterior probability that a scRNA-seq data point was generated  
11 from each cell in the tissue sample and confirmed that the majority of scRNA-seq data points (82.7%) was  
12 assigned to three cells in the tissue with high confidence (>0.8) (**Fig. 2c**, **Supplementary Fig. 1e**). Moreover,  
13 we showed that the scRNA-seq data points were specifically assigned to cells in a small region (a few cell  
14 diameters) of the tissue (**Fig. 2d**, **Supplementary Fig. 1f and g**). These results demonstrated that Perler  
15 accurately reconstructed the spatial expression profiles of the landmark genes using ISH data of all the  
16 landmark genes for training (**Fig. 2e and f**, **Supplementary Fig. 2**) and was capable of doing this via simple  
17 linear mapping.  
18

### 19 **Predictive ability of Perler**

20 We then evaluated the predictive performance of Perler by conducting leave-one-gene-out cross-validation  
21 (LOOCV) in order to confirm whether gene expressions can be predicted following removal of the landmark  
22 gene of interest from the ISH data prior to training (**Supplementary Figs. 3 and 4**). The predictive accuracy  
23 of Perler (aCC = 0.58) was higher than that of Seurat (v.3; aCC = 0.55), Liger (aCC = 0.53), and DistMap  
24 (aCC = 0.44) (**Supplementary Fig. 3**). Moreover, the predictive accuracy can be further improved by  
25 introducing maximum a posteriori (MAP) estimation at the M step in the EM algorithm (see Methods).

26 In addition to the landmark genes, Perler successfully predicted the spatial expression profiles of  
27 non-landmark genes along both anterior–posterior (A–P) and dorsoventral (D–V) axes (**Fig. 3a and b**).  
28 Furthermore, we evaluated the predicted spatial profile of 308 spatially restricted genes (SRGs) proposed by  
29 Bageritz *et al.*<sup>21</sup> (**Supplementary Figs. 4 and 5**) and found that Perler was able to uncover the unknown spatial  
30 gene-expression pattern. Notably, we observed that spatial patterns predicted by Seurat (v.3), Liger, and  
31 DistMap were incomplete. For example, the predicted stripes disappeared in the ventral part of embryos (e.g.,  
32 *abd-A* and *Ubx* in **Fig. 3a**), whereas this issue was not observed with Perler, which accurately predicted the  
33 stripe pattern, even in the ventral part of embryos.  
34  
35



1

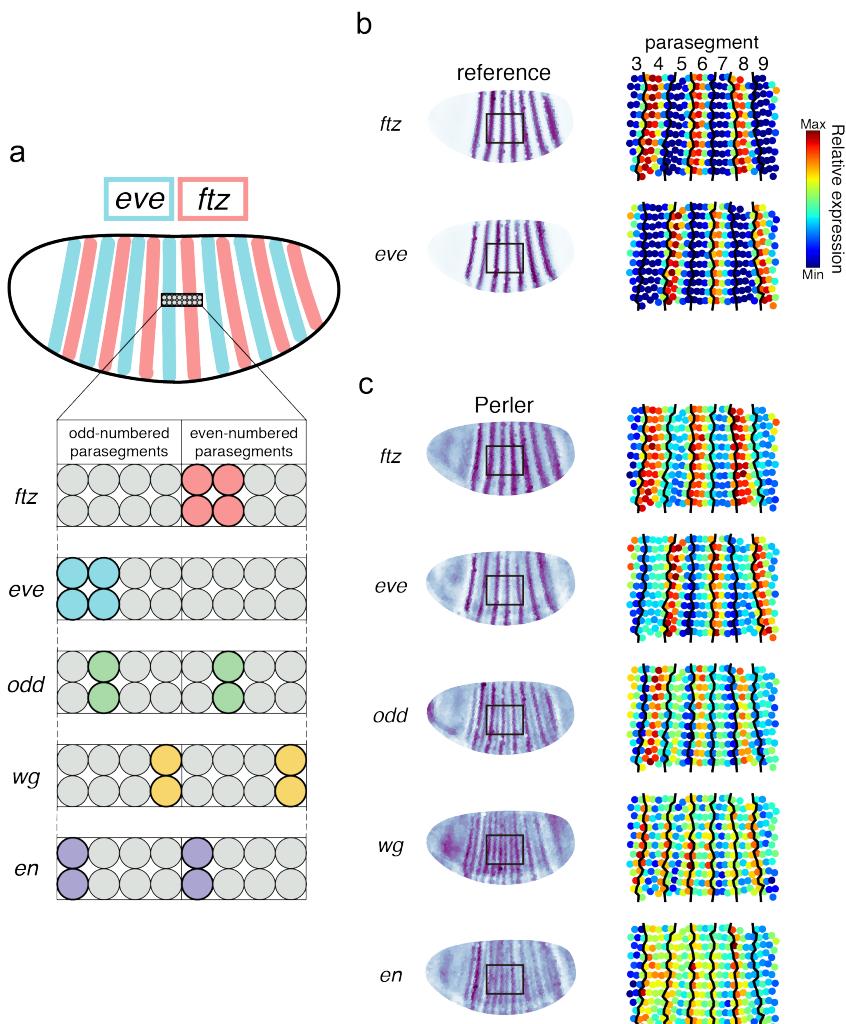
## 2 **Figure 3: Spatial prediction of non-landmark genes**

3 Predictions of non-landmark gene expression showing (a) spatial expression along the A–P and D–V axes and (b) a 14-  
4 stripe pattern according to Perler, Liger, Seurat (v.3), and DistMap. ISH image of *hh* reprinted from Tabata *et al.*<sup>22</sup> under a  
5 Creative Commons License (Attribution: Non-Commercial 4.0 International License).

6

1 **Prediction of 14-stripe patterns of segment-polarity genes**

2 We then presented the spatial predictions of ‘segment-polarity’ genes, which are expressed in a 14-stripe  
3 pattern consistent with the parasegments that subdivide the trunk (main body) region of embryos (**Fig. 3b**)<sup>22–</sup>  
4 <sup>26</sup>. Although the BDTNP reference does not contain information concerning the genes expressed in the 14-  
5 stripe pattern, we found that Perler accurately predicted the spatial expression patterns of these segment-  
6 polarity genes, including *engrailed* (*en*), *wingless* (*wg*), *hedgehog* (*hh*), and *midline* (*mid*) (**Fig. 3b**)<sup>22–26</sup>. By  
7 contrast, all of the previous methods exhibited issues regarding prediction of the 14-stripe patterns. The  
8 predicted patterns demonstrated that DistMap and Seurat (v.3) were unable to predict any 14-stripe patterns,  
9 and that Liger partially predicted 14-stripe patterns, although the ventral part of each stripe was missing (**Fig.**  
10 **3b**). These results suggested that Perler more accurately revealed the spatial gene-expression patterns of non-  
11 landmark genes.



12

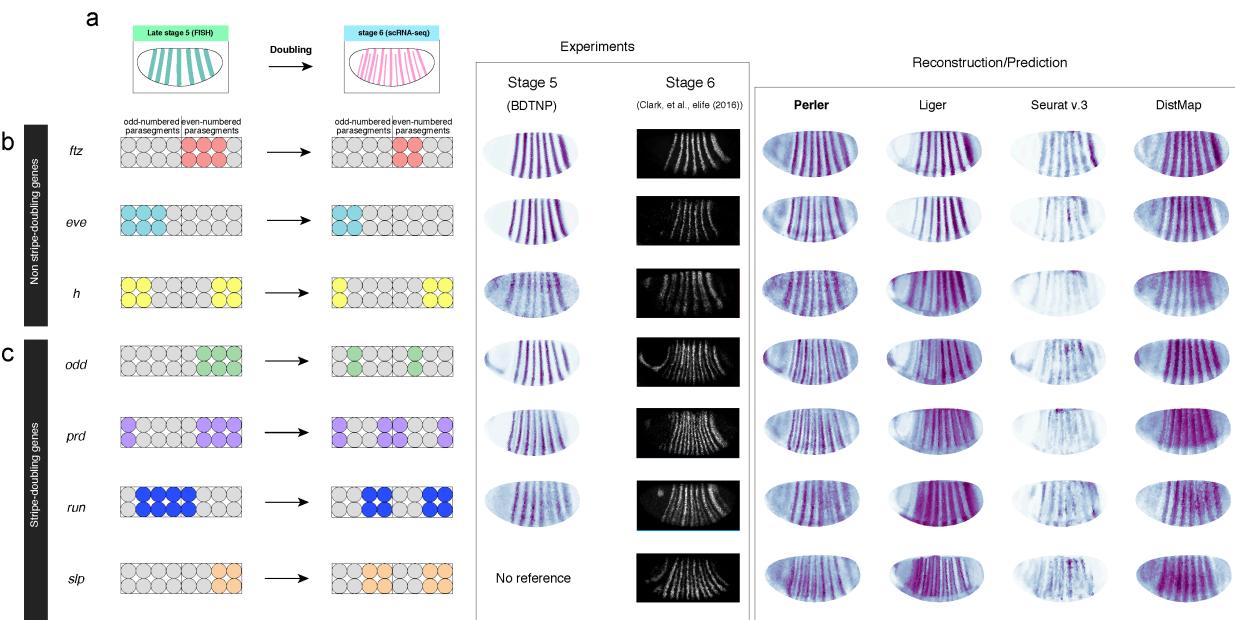
13 **Figure 4: Perler prediction at single-cell resolution**

14 (a) Spatial expression profiles of pair-rule and segment-polarity genes at the single-cell level within each parasegment<sup>27</sup>.  
15 (b) The referenced stripe patterns of *ftz* and *eve*. (c) The reconstructed stripe patterns of *ftz*, *eve*, *odd*, *wg*, and *en* at single-  
16 cell resolution. The left panel shows the spatial gene-expression profiles generated by Perler. The right panel shows the  
17 expanded image of the left panel. The black square in the left panels is the region of interest in the right panels. (b) Black  
18 lines in the right panels indicate the boundaries of each parasegment, which were determined by expression patterns of  
19 *ftz* and *eve* in the reference ISH dataset.

We further analyzed the details of the gene-expression profiles of the segment-polarity genes within each parasegment. Each parasegment shows a four-cell width and is delimited by periodic expression of pair-rule genes and segment-polarity genes at the single-cell width resolution at stage 6<sup>27,28</sup> (**Fig. 4a**). First, we confirmed that the reconstructed patterns of *ftz*, *eve*, and *odd* were consistent with experimental results (**Fig. 4b**). Additionally, the predicted stripes of *wg* were identified adjacent to the predicted stripes of *en*, and the predicted stripes of *en* were identified adjacent to the reconstructed stripes of *odd* (**Fig. 4b and c**). These results were consistent with experimental results<sup>27,28</sup>, strongly supporting the ability of Perler to reveal differences in spatial gene expression at single-cell resolution.

#### Preservation of timing information of scRNA-seq data

We then investigated the effect of timing differences between scRNA-seq (stage 6) and FISH (stage 5) experiments. Although most gene-expression patterns at stage 6 are the same as those at stage 5, several “pair-rule” genes (*odd*, *prd*, *spl*, and *run*) exhibit stripe-doubling from the 7- to the 14-stripe expression patterns during stages 5 and 6<sup>27</sup> (**Fig. 5a**). Accordingly, the scRNA-seq data should intrinsically contain information for the 14-stripe expression pattern. Therefore, we determined whether Perler could reconstruct the 14-stripe pattern from the stage 6 scRNA-seq data.



**Figure 5: Generalization to spatial reference maps by Perler**

(a) Stripe-doubling of pair-rule genes from *Drosophila* developmental stage 5 (FISH experiment) to stage 6 (scRNA-seq experiment). (b, c) Left panels: images of expression changes in (b) non-stripe-doubling genes and (c) stripe-doubling genes. Middle panels: experimental ISH data of stage 5 (BDTNP) and stage 6 (Clark et al.<sup>27</sup>). Right panels: gene expression reconstructed/predicted by Perler, Liger, Seurat (v.3), and DistMap. Note that *spl* expression was predicted. Experimental ISH data were reprinted from Clark et al.<sup>26</sup> under a Creative Commons License (Attribution 4.0 International License).

In our reconstruction, *ftz*, *eve*, and *h* showed a 7-stripe pattern, which was consistent with the previous report<sup>26</sup> showing that these genes do not exhibit stripe-doubling during stages 5 and 6 (**Fig. 5b**). For *odd*, *prd*,

1 and *slp1*, which exhibit stripe-doubling, Perler reconstructions resulted in 14-stripe patterns (**Fig. 5c**).  
2 Additionally, reconstruction of *run* resulted in a partial stripe-doubling pattern, where the third stripe from the  
3 posterior of the embryo was split into two stripes (**Fig. 5c**), surprisingly suggesting that Perler detected the  
4 ongoing phase of a 7-stripe to 14-stripe pattern. These results showed that Perler was able to reconstruct  
5 embryos according to the timing of the scRNA-seq experiment. By contrast, Seurat (v.3) and DistMap  
6 reconstructed every pair-rule gene as 7-stripe patterns (**Fig. 5b and c**). Moreover, Liger reconstructed *odd*, *prd*,  
7 and *slp1* as broad primary seven stripes with weak secondary seven stripes, which were so obscure that it was  
8 difficult to distinguish 14 stripes, and reconstructed *run* as a 7-stripe pattern (**Fig. 5c**). These results indicated  
9 that previous methods reconstructed embryos according to the timing of FISH experiments rather than that of  
10 scRNAseq experiments. Taken together, these findings showed that Perler successfully reconstructed spatial  
11 gene-expression profiles according to the timing of scRNA-seq experiments (stage 6), regardless of the timing  
12 of FISH experiments (stage 5), while all other methods reconstructed those at the timing of FISH experiments.  
13 We concluded that Perler has the ability to not over-fit to ISH data and robustly preserve timing information  
14 in scRNA-seq data.

15

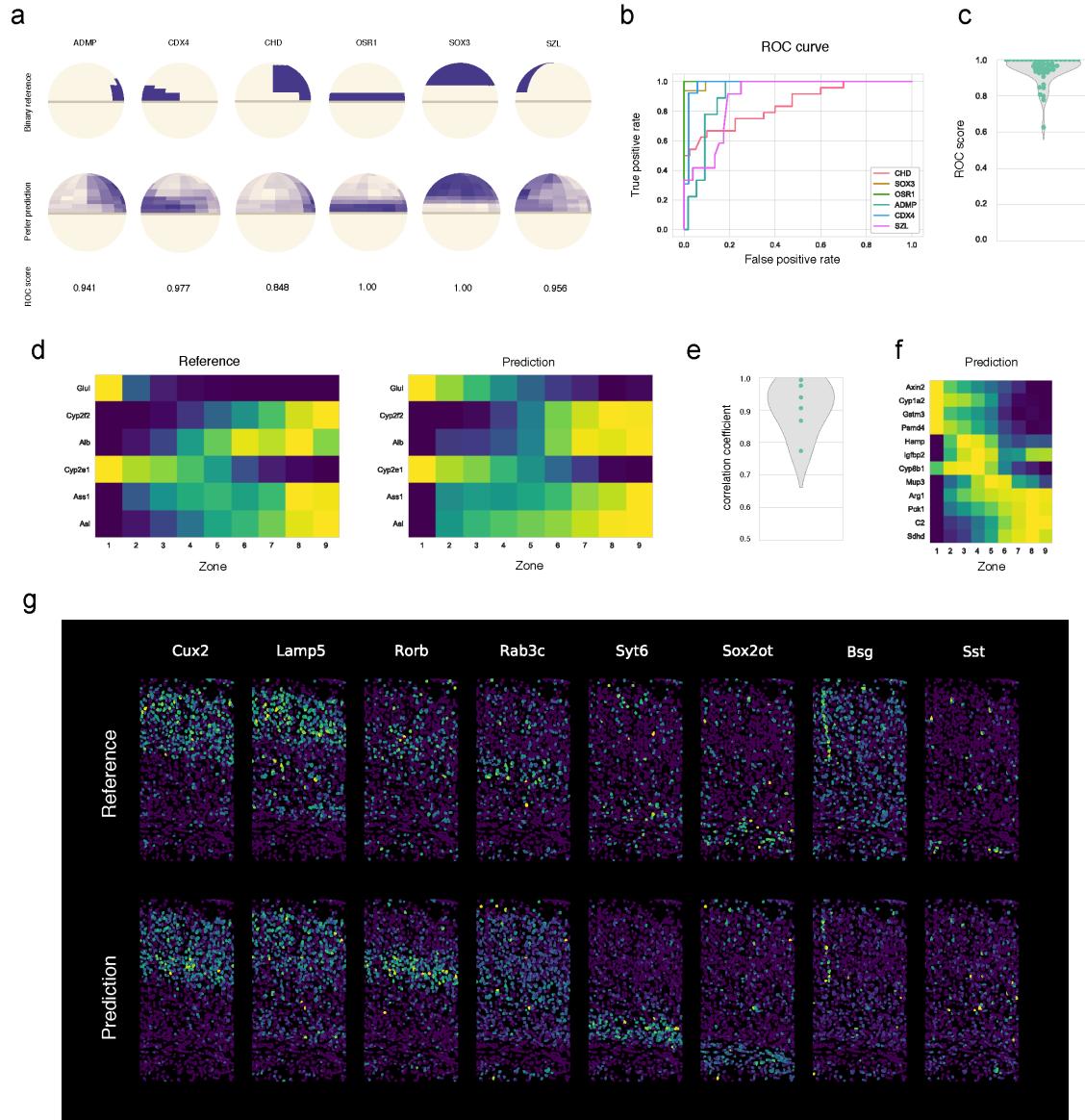
### 16 *Application to other datasets*

17 To evaluate Perler applicability to other datasets, we evaluated it using three published datasets. First, we  
18 applied Perler to the zebrafish embryo datasets, in which the spatial reference map was binarized based on  
19 traditional measurement by ISH<sup>4</sup> (**Fig. 6a**; see Methods). Cross-validation demonstrated that Perler accurately  
20 predicted spatial gene-expression profiles compatible with Seurat (v.1)<sup>4</sup> [median receiver operating  
21 characteristic (ROC) score = 0.96], even using the binary spatial reference map (**Fig. 6a–c**).

22 We then applied Perler to mammalian liver datasets, in which the spatial reference map was measured  
23 by single-molecule (sm)FISH<sup>7</sup> (**Fig. 6d**, see Methods). Cross-validation showed that the predictive accuracy  
24 (aCC = 0.91) was sufficiently high (**Fig. 6e**), and that Perler successfully predicted both monotonic and non-  
25 monotonic gene-expression gradients (**Fig. 6f**)<sup>7</sup>. Finally, we applied Perler to adult mouse visual cortex  
26 datasets, in which the single-cell resolution ISH data for 1,020 genes was measured by recent *in situ* technology  
27 (STARmap<sup>15</sup>), and scRNA-seq data for 14,739 cells available from the Allen Brain Atlas<sup>16</sup>. Cross-validation  
28 revealed that Perler predicted the spatial expression patterns of genes according to both layer-specific  
29 expression and cell-type-specific expression in brain cortex (**Fig. 6g**). These results suggested that Perler is  
30 applicable to prediction using high-dimensional spatial reference maps.

31 Taken together, the findings support Perler as a powerful tool for predicting spatial gene-expression  
32 profiles in any multicellular system with general applicability to any type of ISH data (e.g., binary or  
33 continuous, low to high dimension, and single-cell to tissue-level resolution).

34



1

## 2 **Figure 6: Applications of PERLER to other data**

3 (a–c) Application of Perler for early zebrafish embryo data. (a) LOOCV experiments. The upper and lower panels show the  
 4 referenced ISH data and predicted gene-expression profiles, respectively. (b) ROC curve for the LOOCV experiments for  
 5 the genes shown in (a). (c) Violin plot for the predictive accuracies of Perler for the LOOCV experiments for all genes in  
 6 the reference ISH data according to ROC score. (d–f) Application of Perler for mammalian liver lobules. (d) LOOCV  
 7 experiments. The left and right panels show the reference ISH data and predicted gene-expression profiles, respectively.  
 8 All genes from the ISH data are shown. (e) Violin plot for the predictive accuracies of Perler for the LOOCV experiments  
 9 for all genes in the reference ISH data. (f) Prediction of non-landmark genes. In addition to the monotonic gene-expression  
 10 profiles, non-monotonic gene-expression profiles are observed (*Hamp*, *Igfbp2*, *Cyp8b1*, and *Mup3*)<sup>7</sup>. (g) Application of  
 11 Perler for the mouse visual cortex. The upper and lower panels show the referenced ISH data and the predicted gene-  
 12 expression profiles, respectively.

13

14

## 1 Discussion

2 In this study, we developed a model-based computational method (Perler) that predicts genome-wide spatial  
3 transcriptomes. Perler sequentially conducted a two-step computation, with the first step mapping ISH data  
4 points to the scRNA-seq space according to the generative linear model by EM algorithm (**Fig. 1b**), and the  
5 second step optimizing the weighting function used to predict spatial transcriptomes according to weighted  
6 scRNA-seq data points (**Fig. 1c and d**). Using a dataset for early *Drosophila* embryos, we demonstrated that  
7 Perler accurately reconstructed and predicted genome-wide spatial transcriptomes and was able to robustly  
8 preserve the timing information in scRNA-seq data. (**Figs. 2–5**). Moreover, we showed that in any multicellular  
9 system, Perler displayed broad applicability to any type of ISH data (**Fig. 6**).  
10

### 11 *Difference from existing methods*

12 Several studies have analyzed spatial gene-expression profiles from scRNA-seq data. Here, we discuss their  
13 differences from Perler (**Supplementary Table 1**). First, Perler does not require binarization of both ISH and  
14 scRNA-seq data, which differs from Seurat (v.1)<sup>4</sup>, the method described by Achim *et al.*<sup>6</sup>, and DistMap<sup>5</sup>, all  
15 of which binarize either or both data sets and result in loss of information related to continuous gene expression.  
16 We showed that Perler significantly improved the predictive accuracy of the *Drosophila* spatial transcriptome  
17 as compared with DistMap (**Fig. 2**, **Supplementary Fig. 3**). Additionally, even with data for zebrafish embryos  
18 only available in a traditional binary ISH dataset, Perler showed similar predictive performance with that of  
19 Seurat (v.1) (**Fig. 6**).

20 Perler does not need ISH data with a gene-expression distribution for each cell or subregion, which  
21 differs from the method of Halpern *et al.*<sup>7</sup>, which repeated smFISH experiments in order to sample distributions  
22 of landmark gene expression on a one-dimensional radial coordinate in hepatic lobules. Using the sampling  
23 distributions, their Bayesian method computed posterior probabilities of each scRNA-seq data point being  
24 generated from cells/subregions in the tissue. However, repeating smFISH experiments is obviously labor  
25 intensive in the case of two- or three-dimensional tissue. Perler addresses this sampling problem by estimating  
26 gene-expression distributions in each cell from the tissue sample using the generative linear model. This  
27 enables Perler to calculate posterior probabilities for two- and three-dimensional tissues (**Figs. 2c and 6**).

28 Perler is a model-based method, which differs from Seurat (v.3)<sup>9</sup> and Liger<sup>10</sup>, both based on model-  
29 free mapping between ISH and scRNA-seq data (e.g., CCA and NMF methods). Their model-free mapping  
30 addresses gene expression as continuous variables with applicability to any kind of multicellular system;  
31 however, these methods freely map ISH data to scRNA-seq data without any assumptions (i.e., they do not  
32 account for latent relationships between the two datasets). Indeed, we showed that Seurat v.3 and Liger over-  
33 fit to the timing of ISH experiments by focusing on the stripe-doubling of pair-rule genes in *Drosophila* (**Fig.**  
34 **5**). To guarantee generalized performance, we introduced generative linear modeling with biologically  
35 interpretable constraints and statistically reasonable distances. For the former, expression levels are linearly  
36 correlated between ISH and scRNA-seq measurements with gene-specific sensitivity, background signals, and  
37 noise intensity. For the latter, the pairwise distances between ISH and scRNaseq data points were evaluated

1 in a variance-scaled manner using Mahalanobis' metric of Gaussian mixture distribution (see Methods).  
2 Consequently, Perler avoided the over-fitting problem encountered by previous methods. Indeed, Perler  
3 reconstructed the stripe-doubling of pair-rule genes according to the timing of scRNA-seq data (**Fig. 5**).

4 It is worth mentioning a recent method called novoSpaRc<sup>29</sup>. This method proposed a new concept  
5 for predicting spatial expression patterns using the physical information of cells in tissue, which enables these  
6 predictions with little or no information regarding ISH gene-expression patterns. However, in practice, their  
7 predictive ability using *Drosophila* scRNA-seq data is unsatisfactory at single-cell resolution; therefore, this  
8 concept of using cellular information remains challenging. As a focus of future study, it would be interesting  
9 to extend our generative model to introduce prior knowledge of physical information.

10

### 11 ***Application to multi-omics analysis***

12 We demonstrated that Perler can integrate two distinct datasets of scRNA-expression profiles while also  
13 avoiding overfitting to the reference. These features suggest that Perler could be a suitable theoretical  
14 framework for integrating not only two RNA-expression datasets but also two single-cell datasets with  
15 different modalities, such as chromatin accessibility measured by a single-cell assay for transposase-accessible  
16 chromatin using sequencing and DNA methylation measured by chromatin immunoprecipitation sequencing.  
17 Particularly in terms of multi-omics analysis, where datasets from two different modalities do not exactly  
18 match and are often sampled from different individuals and using different time intervals<sup>30,31</sup>, Perler can  
19 potentially help integrate different types of single-cell genomics data. Thus, Perler provides a powerful and  
20 generalized framework for revealing the heterogeneity of multicellular systems.

21

22

## 1 Methods

2 We developed a novel method to reconstruct spatial gene-expression profiles from an scRNA-seq dataset via  
3 comparison with a spatial reference map measured by ISH-based methods. In the spatial reference map,  
4 landmark gene-expression vectors ( $D$  genes; e.g.,  $D = 84$  in early *D. melanogaster* embryos) are available for  
5 all cells, whose locations in the tissue are known. The landmark gene-expression vector of cell  $k$  is represented  
6 as  $\mathbf{h}_k=(h_{k,1}, h_{k,2}, \dots, h_{k,D})^T$ , where cells are indexed by  $k$  ( $k \in \{1, 2, \dots, K\}$ ), and  $K$  is the total number of cells in  
7 the tissue of interest. By contrast, in an scRNA-seq dataset, genome-wide expression ( $D'$  genes; e.g.,  $D' =$   
8 8924 in early *D. melanogaster* embryos) lack information regarding cell location in tissue. The genome-wide  
9 expression vector of cell  $n$  is represented as  $\mathbf{y}_n=(y_{n,1}, y_{n,2}, \dots, y_{n,D'})^T$ , where cells are indexed by  $n$  ( $n \in \{1, 2, \dots,$   
10  $N\}$ ), and  $N$  is the total number of cells used for scRNA-seq measurement.

11

### 12 Observation model

13 We modeled the difference between scRNA-seq and ISH measurements as

$$14 \quad y_i = a_i h_i + b_i + c_i \xi_i \quad (1)$$

15 where  $y_i$  and  $h_i$  indicate expression levels of landmark gene  $i$  measured by scRNA-seq and ISH experiments,  
16 respectively;  $\xi_i$  indicates Gaussian noise with zero mean and unit variance; and  $a_i$ ,  $b_i$ , and  $c_i$  are constant  
17 parameters for gene  $i$ , which are interpreted as scale difference amplification rates, background signals, and  
18 noise intensities, respectively.

19 We reduced the dimensionality of the genes to change equation (1) to

$$20 \quad x_j = a_j r_j + b_j + c_j \xi_j \quad (2)$$

21 where  $x_j$  and  $r_j$  indicate expression levels of metagene  $j$  for scRNA-seq and ISH in the lower dimensional space,  
22  $j \in \{1, 2, \dots, M\}$ ; and  $M$  indicates the number of metagenes. In vector-matrix representation, the equation (2)  
23 is written as

$$24 \quad \mathbf{x} = \mathbf{A}\mathbf{r} + \mathbf{b} + \mathbf{C}\xi \quad (3)$$

25 where  $\mathbf{x} = (x_1, x_2, \dots, x_M)^T$ ,  $\mathbf{r} = (r_1, r_2, \dots, r_M)^T$ ,  $\mathbf{A} = \text{diag}(a_1, a_2, \dots, a_M)$ ,  $\mathbf{b} = (b_1, b_2, \dots, b_M)^T$ ,  $\mathbf{C} = \text{diag}(c_1, c_2,$   
26  $\dots, c_M)$ , and  $\xi = (\xi_1, \xi_2, \dots, \xi_M)^T$ .

27

### 28 Metagene representation in lower dimensional space

29 The dimensionalities of both scRNA-seq and reference data were reduced by PLSC analysis<sup>18</sup>. PLSC can  
30 extract the correlated coordinates from both datasets. In PLSC analysis, the cross-correlation matrix of scRNA-  
31 seq and ISH data is first calculated as

$$32 \quad \mathbf{W} = \mathbf{Y}^T \mathbf{H} \quad (4)$$

33 where  $\mathbf{Y}$  and  $\mathbf{H}$  indicate a  $D \times N$  scRNA-seq data matrix with  $D$  landmark genes and  $N$  cells, and a  $D \times K$  ISH  
34 data matrix with  $D$  landmark genes and  $K$  cells, respectively.  $\mathbf{W}$  is then subjected to singular value  
35 decomposition as

$$36 \quad \mathbf{W} \simeq \mathbf{U}^T \Delta \mathbf{V} \quad (5)$$

37

1 where  $\mathbf{U}$ ,  $\Delta$ , and  $\mathbf{V}$  indicate the  $M \times N$  singular vector matrices, the  $M \times M$  diagonal matrix, and  $M \times K$  singular  
2 vector matrices, respectively, with  $M$  representing the reduced dimension (i.e., the number of metagenes). In  
3 this study, the metagene vectors for scRNA-seq ( $\mathbf{x}_n$ ) and the reference data ( $\mathbf{r}_k$ ) were respectively calculated  
4 by

5 
$$\mathbf{x}_n = \Delta \mathbf{u}_n \quad (6)$$

6 
$$\mathbf{r}_k = \Delta \mathbf{v}_k \quad (7)$$

7 where  $\mathbf{u}_n$  and  $\mathbf{v}_k$  indicate the  $n^{\text{th}}$  row vector of  $\mathbf{U}$  and the  $k^{\text{th}}$  row vector of  $\mathbf{V}$ , respectively.  
8

9 ***A Gaussian mixture model (GMM) for scRNA-seq observation***

10 We used equation (3) to transform ISH observations into scRNA-seq observations. To infer from which cells  
11 in the tissue the scRNA-seq observations originated, we developed a generative model for metagene-  
12 expression vectors for scRNA-seq data  $\mathbf{x}$ , which was expressed by a K-components GMM:

13 
$$P(\mathbf{x}|\theta) = \sum_{k=1}^K \pi_k N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}), \quad (8)$$

14 where

15 
$$\boldsymbol{\mu}_k = \mathbf{A}\mathbf{r}_k + \mathbf{b}, \quad (9)$$

16  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2)$  ( $\sigma_j=c_j$ ),  $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  indicates a multivariate Gaussian distribution with mean and  
17 variance-covariance matrix  $\boldsymbol{\Sigma}$ , and  $\pi_k$  is the probability that  $\mathbf{x}$  originated from cell  $k$  in the tissue. Note that  $\mathbf{A}$ ,  
18  $\mathbf{b}$ , and  $\boldsymbol{\Sigma}$  are unknown parameters that need to be estimated.

19 The log of likelihood function of this GMM model is given by

20 
$$L(\boldsymbol{\theta}) = \sum_{n=1}^N \ln \left( \sum_{k=1}^K \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}) \right) \quad (10)$$

21 where  $\boldsymbol{\theta}$  indicates a set of the parameters  $\boldsymbol{\theta} \in \{\boldsymbol{\pi}, \mathbf{A}, \mathbf{b}, \boldsymbol{\Sigma}\}$  and  $\boldsymbol{\pi}=(\pi_1, \pi_2, \dots, \pi_M)^T$ .  
22

23 ***EM algorithm (the first step in Perler)***

24 To estimate the unknown parameters ( $\boldsymbol{\pi}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\boldsymbol{\Sigma}$ ), we maximize the log likelihood function using the EM  
25 algorithm. In the E step, based on the current parameter values, we calculated the responsibility, which  
26 represents the posterior probability that scRNA-seq vector  $\mathbf{x}_n$  was derived from cell  $k$  in the tissue as  
27

28 
$$\gamma_{nk} = \frac{\pi_k N(\mathbf{x}_n|\mathbf{A}^{(\text{old})}\mathbf{r}_k + \mathbf{b}^{(\text{old})}, \boldsymbol{\Sigma}^{(\text{old})})}{\sum_j^K \pi_j N(\mathbf{x}_n|\mathbf{A}^{(\text{old})}\mathbf{r}_k + \mathbf{b}^{(\text{old})}, \boldsymbol{\Sigma}^{(\text{old})})}. \quad (11)$$

29  
30 In the M step, we optimize the parameter values in order to maximize the log likelihood function based on the  
31 current responsibilities. These parameter values are updated as follows:  
32

$$\pi_k^{(\text{new})} = \frac{\sum_j^K \gamma_{nk}}{N}, \quad (12)$$

1 
$$a_i^{(new)} = \frac{\psi_i - \chi_i b_i^{(new)}}{\omega_i}, \quad (13)$$

2 
$$b_i^{(new)} = \frac{\omega_i \phi_i - \psi_i \chi_i}{N \omega_i - \chi_i^2}, \quad (14)$$

3 
$$\sigma_i^{2(new)} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left( x_{ni} - a_i^{(new)} r_{ki} - b_i^{(new)} \right)^2, \quad (15)$$

4 where

5 
$$\phi_i = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} x_{ni}, \quad (16)$$

6 
$$\chi_i = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} r_{ni}, \quad (17)$$

7 
$$\psi_i = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} x_{ni} r_{ni}, \quad (18)$$

8 
$$\omega_i = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} r_{ni}^2. \quad (19)$$

9 The detailed derivation for these equations is presented in a later subsection. The E and M steps iterate until  
10 the log likelihood function converges, after which the obtained estimated parameters  $\hat{\theta} \in \{\hat{\pi}, \hat{\mathbf{A}}, \hat{\mathbf{b}}, \hat{\Sigma}\}$ .  $\hat{\mathbf{u}}_k$  are  
11 given as

12 
$$\hat{\mathbf{u}}_k = \hat{\mathbf{A}} \mathbf{r}_k + \hat{\mathbf{b}}, \quad (20)$$

13 describing the mapped metagene-expression vector of cell  $k$  measured by ISH. Note that  $\hat{\mathbf{u}}_k$  is the metagene-  
14 expression vector in the scRNA-seq space.

15

### 16 ***Spatial reconstruction (the second step in Perler)***

17 We reconstructed/predicted the gene-expression vector by weighted averaging all scRNA-seq data points as

18 
$$\bar{\mathbf{y}}_k = \sum_{n=1}^N \frac{w_{nk} \mathbf{y}_n}{\sum_{j=1}^N w_{jk}}, \quad (21)$$

19 where  $\mathbf{y}_n$  indicates the  $n^{\text{th}}$  scRNA-seq data point ( $D$ -component vector).  $w_{nk}$  is calculated by  
20

21 
$$w_{nk} = \frac{\pi_k \exp(-\alpha D_{nk}^2 - \beta D_{nk} - \delta)}{\sum_{j=1}^K \pi_j \exp(-\alpha D_{nj}^2 - \beta D_{nj} - \delta)}, \quad (22)$$

22 where  $\alpha$ ,  $\beta$ , and  $\delta$  are positive constants. Note that  $\delta$  in the numerator and denominator of equation (22) are  
23 canceled out.  $D_{nk}$  indicates Mahalanobis' distance between scRNA-seq data point  $\mathbf{x}_n$  and cell  $k$ :

1                    $D_{nk} = \sqrt{(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)}.$                    (23)

2  
3     If  $\alpha = 1/2$  and  $\beta = 0$ ,  $w_{nk}$  is exactly the posterior probability that scRNA-seq data point  $\mathbf{x}_n$  is generated by cell  
4      $k$ . Note that equation (21) has a similar structure to the Nadaraya–Watson model<sup>13</sup>. Values of  $\alpha$  and  $\beta$  are  
5     determined by cross-validation.

6  
7     **Hyperparameter optimization**  
8     We optimized the hyperparameters  $\alpha$  and  $\beta$  of the weighting function by LOOCV in order to fit the predicted  
9     gene expression to the referenced gene expression measured by ISH. To this end, we removed one of the  
10    landmark genes from the ISH data and used this dataset to predict the spatial gene-expression profile of the  
11    removed landmark gene with the fixed hyperparameters in Perler. This LOO prediction was repeated for every  
12    landmark gene. We then quantitatively evaluated the predictive performance of these hyperparameters  
13    according to the mutual information existing between the predicted expression and referenced expression of  
14    all landmark genes:

15                    $J = -\frac{1}{2} \sum_i^D \ln\{1 - \rho_i(\alpha, \beta)^2\},$                    (24)

16     where  $J$  is the approximated mutual information between the predicted and referenced gene expression.  $\rho_i(\alpha,$   
17      $\beta)$  indicates the Pearson's correlation coefficient between the predicted spatial expression pattern of each  
18     landmark gene  $i$  and its reference ISH data as

19                    $\rho_i(\alpha, \beta) = \frac{\sum_k^K (\bar{y}_{ki} - \langle \bar{y}_i \rangle)(h_{ki} - \langle h_i \rangle)}{\sqrt{\sum_k^K (\bar{y}_{ki} - \langle \bar{y}_i \rangle)^2} \sqrt{\sum_k^K (h_{ki} - \langle h_i \rangle)^2}},$                    (25)

20     where

21                    $\langle \bar{y}_i \rangle = \frac{1}{K} \sum_{k=1}^K \bar{y}_{ki}, \text{ and}$                    (26)

22                    $\langle h_i \rangle = \frac{1}{K} \sum_{k=1}^K h_{ki}.$                    (27)

23     The derivation of  $J$  is described in a later subsection. Here, we optimized  $\alpha$  and  $\beta$  by grid search in order to  
24     maximize the mutual information,  $J$ . We then used the optimized hyperparameters to predict the spatial profile  
25     of non-landmark genes (**Fig. 3**, **Supplementary Fig. 5**). To evaluate the predictive performance of Perler (**Fig.**  
26     **3**), we removed each landmark gene from the mutual information and re-optimized the hyperparameters. This  
27     re-optimization is repeated for every landmark gene. Note that for the zebrafish embryo data, we used the ROC  
28     score instead of the correlation coefficient, because only the binary ISH data was available. Additionally, for

1 the mouse visual cortex data, we conducted 10-fold cross-validation because of the massive computational  
2 cost of LOOCV for the large number of landmark genes (1,020 genes).

3

4 **Data acquisition and preprocessing**

5 For *D. melanogaster* reconstruction, we used scRNA-seq and ISH data at *Drosophila* Virtual Expression  
6 eXplorer (DVEX) (<https://shiny.mdc-berlin.de/DVEX/><sup>5</sup>), which was originally used for DistMap<sup>5</sup>. In these  
7 data sets, the number of scRNASeq data points is 1297, whereas the number of cells to be estimated in the  
8 embryos is 3039. The expressed mRNA counts in this scRNA-seq dataset were already log normalized  
9 according to the total number of unique molecular identifiers (UMIs) for each cell. For each gene, we  
10 subtracted the average expression from the scRNA-seq data. Additionally, the ISH data were log-scaled and  
11 subtracted average expression from this ISH data, as same as the scRNA-seq data.

12 For reconstruction of the early zebrafish embryos, we acquired the public scRNA-seq and ISH data  
13 from the Satija Lab homepage (<https://satijalab.org/><sup>4</sup>), with these data originally used by Seurat (v.1)<sup>4</sup>. In these  
14 data, the number of scRNASeq data points is 851, whereas the number of subregions to be estimated in the  
15 embryos is 64. Note that the ISH data were binary. Similar to the *Drosophila* data, we log-scaled both scRNA-  
16 seq and ISH datasets and subtracted the average expression of each gene.

17 For reconstruction of the mammalian liver, we used scRNA-seq and smFISH data provided by  
18 Halpern *et al.*<sup>7</sup>. In these data, the number of scRNASeq data points is 1415, whereas the number of zones to be  
19 estimated in the embryos is 9. Because multiple samples were provided in the smFISH data, we calculated  
20 their average at each tissue location for Perler, followed by log-scaling both the scRNA-seq and smFISH data  
21 and subtracting the average expression of each gene.

22 For reconstruction of the mouse visual cortex, we used scRNA-seq data provided by the Allen Brain  
23 Institute<sup>16</sup> and smFISH data provided by Wang *et al.*<sup>15</sup>, respectively, which were originally used for Seurat  
24 (v.3)<sup>9</sup>. The number of scRNASeq data points is 14739, whereas the number of cells to be estimated in the  
25 cortex is 1549. We log-scaled both the scRNA-seq and smFISH data and subtracted the average expression of  
26 each gene.

27

28 **Data visualization**

29 For *D. melanogaster*, we visualized the reconstructed gene-expression profile at single-cell resolution by using  
30 the three-dimensional coordinates of all cells from DVEX (<https://shiny.mdc-berlin.de/DVEX/><sup>5</sup>). Because the  
31 embryo is bilaterally symmetric, we mapped the reconstructed spatial gene-expression levels of the 3,039 cells  
32 in the right-half embryo. According to the previous study<sup>5</sup>, we then mirrored the spatial gene-expression levels  
33 of the right-half cells to the remaining cells in left-half embryo. In the case of the early zebrafish embryos, we  
34 visualized the reconstructed gene expression using the ‘zf.insitu.vec.lateral’ function of Seurat (v.  $\geq 1.2$ )<sup>4</sup>. In  
35 the case of the mammalian liver, we visualized the reconstructed gene expression as a heatmap. In the case of  
36 the mouse visual cortex, we visualized the reconstructed gene expression at single-cell resolution. We used  
37 two-dimensional coordinates of all cells within cortical slices provided by Wang *et al.*<sup>15</sup>.

38

1    **Derivation of the EM algorithm**

2    The goal of the EM algorithm is to maximize the likelihood function  $p(\mathbf{X}|\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ , where  $\mathbf{X} = \{\mathbf{x}_1,$   
3     $\mathbf{x}_2, \dots, \mathbf{x}_N\}$  and  $\boldsymbol{\theta} = \{\boldsymbol{\pi}, \mathbf{A}, \mathbf{b}, \boldsymbol{\Sigma}\}$ . The generative model of scRNA-seq data point  $\mathbf{x}$  with latent variables  $\mathbf{z}$  is  
4    formulated, as follows. The probability distribution of  $\mathbf{z}$  is

5                  
$$P(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}, \quad (28)$$

6    where  $\mathbf{z}$  is a vector in a one-of- $K$  representation that shows from which cells/regions in tissue a scRNA-seq  
7    sample originated;  $z_k$  is the  $k^{\text{th}}$  element of  $\mathbf{z}$ ;  $K$  is the number of the elements of the latent variables  $\mathbf{z}$  equal to  
8    the number of cells in the tissue; and  $\pi_k$  is probability that  $z_k = 1$ . The probability distribution of  $\mathbf{x}$  conditioned  
9    by  $\mathbf{z}$  is

10                 
$$P(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K N(\mathbf{x}|\mathbf{A}\mathbf{r}_k + \mathbf{b}_k, \boldsymbol{\Sigma})^{z_k}, \quad (29)$$

11    where  $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  indicates a Gaussian distribution with mean  $\boldsymbol{\mu}$  and variance  $\boldsymbol{\Sigma}$ ;  $\mathbf{A}_k$  is the  $M \times M$  diagonal matrix;  
12    $\mathbf{b}_k$  indicates the  $M$  elements vector in equation (3); and  $\mathbf{r}_k$  indicates the  $M$  elements vector describing the  
13   metagene-expression level in cell  $k$ . The joint probability distribution of  $\mathbf{x}$  and  $\mathbf{z}$  is

14                 
$$P(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K \{\pi_k N(\mathbf{x}|\mathbf{A}\mathbf{r}_k + \mathbf{b}_k, \boldsymbol{\Sigma})\}^{z_k}. \quad (30)$$

15    Note that the marginalized distribution of  $\mathbf{z}$  becomes equation (8). The likelihood function for the complete  
16   dataset  $\{\mathbf{X}, \mathbf{Z}\}$  is given as

17                 
$$P(\mathbf{X}, \mathbf{Z}) = \prod_{n=1}^N \prod_{k=1}^K \{\pi_k N(\mathbf{x}_n|\mathbf{A}\mathbf{r}_k + \mathbf{b}, \boldsymbol{\Sigma})\}^{z_{nk}}, \quad (31)$$

18    where  $\mathbf{Z} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N\}$ . Therefore, the expectation of its log-likelihood function over the posterior  
19   distribution of  $P(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(old)})$  becomes

20                 
$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(old)}) = \sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{X}, \boldsymbol{\theta}^{(old)}) \ln P(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \quad (32)$$

21                 
$$= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \ln \pi_k + \gamma_{nk} \ln N(\mathbf{x}_n|\mathbf{A}\mathbf{r}_k + \mathbf{b}, \boldsymbol{\Sigma}), \quad (33)$$

22    where  $\gamma_{nk}$  is the expectation of  $z_{nk}$  over  $P(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(old)})$  given as

23                 
$$\gamma_{nk} = \sum_{\mathbf{z}_n} z_{nk} P(\mathbf{z}_n|\mathbf{x}_n, \boldsymbol{\theta}^{(old)}). \quad (34)$$

24    According to Bayes' theorem,

25                 
$$P(\mathbf{z}_n|\mathbf{x}_n, \boldsymbol{\theta}^{(old)}) = \frac{P(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta}^{(old)}) P(\mathbf{z}_n, \boldsymbol{\theta}^{(old)})}{\sum_{\mathbf{z}_n} P(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta}^{(old)}) P(\mathbf{z}_n, \boldsymbol{\theta}^{(old)})}, \quad (35)$$

1 where  $P(z_{nk}=1|\mathbf{x}_n)$  becomes equation (11).

2 In the E step,  $\gamma_{nk}$  is calculated based on the current parameter values of  $\boldsymbol{\theta}^{(old)}$ . In the M step, we update  
3 the parameter values  $\boldsymbol{\theta}$  by maximizing the Q-function as

4 
$$\boldsymbol{\theta}^{(new)} = \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(old)}). \quad (36)$$

5 The maximization of  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(old)})$  with respective to  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\boldsymbol{\Sigma}$  is achieved by  $\partial Q/\partial \mathbf{A} = 0$ ,  $\partial Q/\partial \mathbf{b} = 0$ , and  
6  $\partial Q/\partial \boldsymbol{\Sigma} = 0$ , leading to equations (13–19).  $\boldsymbol{\pi}$  is updated by introducing a Lagrange multiplier to enforce the  
7 constraint  $\sum_{k=1}^K \pi_k = 1$ , leading to equation (12).

8

9 **Derivation of mutual information**

10 We derived equation (24) by approximating the following mutual information between the reconstructed  
11 spatial expression pattern of the landmark genes and their reference map:

12 
$$I(\bar{\mathbf{y}}, \mathbf{h}) = \int \int P(\bar{\mathbf{y}}, \mathbf{h}) \ln \frac{P(\bar{\mathbf{y}}, \mathbf{h})}{P(\bar{\mathbf{y}})P(\mathbf{h})} d\bar{\mathbf{y}} d\mathbf{h}, \quad (37)$$

13 where  $\bar{\mathbf{y}} = (\bar{y}_1, \bar{y}_2, \dots, \bar{y}_D)$ ,  $\mathbf{h} = (h_1, h_2, \dots, h_D)$ , and  $\bar{y}_i$  and  $h_i$  indicate random variables representing the  
14 predicted and referenced expression levels of landmark gene  $i$ , respectively, and  $P(\bar{\mathbf{y}}, \mathbf{h})$  indicates the joint  
15 probability distribution of  $\bar{\mathbf{y}}$  and  $\mathbf{h}$ . Here, we assumed that spatial expressions of landmark genes are  
16 independent from one another, which leads to

17 
$$I(\bar{\mathbf{y}}, \mathbf{h}) = \sum_i^D \int \int P(\bar{y}_i, h_i) \ln \frac{P(\bar{y}_i, h_i)}{P(\bar{y}_i)P(h_i)} d\bar{y}_i dh_i. \quad (38)$$

18 We calculated  $I(\bar{\mathbf{y}}, \mathbf{h})$  by assuming  $P(\bar{y}_i, h_i)$  as a bivariate Gaussian distribution and obtained

19 
$$I(\bar{\mathbf{y}}, \mathbf{h}) = -\frac{1}{2} \sum_i^D \ln(1 - \rho_i(\alpha, \beta)^2), \quad (39)$$

20 where  $\rho_i(\alpha, \beta)$  denotes the calculated Pearson's correlation coefficient calculated.

21

22

## 1 **Code availability**

2 The Python code of Perler will be available when published.

3

## 4 **Acknowledgements**

5 This study was supported in part by the Cooperative Study Program of Exploratory Research Center on Life and  
6 Living Systems (ExCELLS) (program Nos.18-201, 19-102, and 19-202 to H.N.); a Grant-in-Aid for Young  
7 Scientists (B) (16K16147 and 19H04776 to H.N.), a Grant-in-Aid for Scientific Research (B) (17KT0021 to T.K.)  
8 and a JSPS research fellowship for young scientist (to S.S.) from the Japan Society for the Promotion of Science  
9 (JSPS); the Naito Foundation (to T.K.); and the Keihanshin Consortium for Fostering the Next Generation of Global  
10 Leaders in Research (K-CONNEX) established by the program of Building of Consortia for the Development of  
11 Human Resources in Science and Technology, MEXT (to T.K.).

12

## 13 **Author Contributions**

14 H.N., S.S. and T.K. conceived the project. Y.O., H.N. and K.N. developed the method, Y.O. implemented the  
15 software, and Y.O. and S.S. analyzed data. Y.O. and H.N. wrote the manuscript with input from all authors.

16

## 17 **Competing Interests**

18 The authors declare no competing interests.

19

20

## 1 References

- 2 1. Gilmour, D., Rembold, M. & Leptin, M. From morphogen to morphogenesis and back. *Nature* **541**,  
3 311–320 (2017).
- 4 2. Halpern, K. B. *et al.* Paired-cell sequencing enables spatial gene expression mapping of liver  
5 endothelial cells. *Nat. Biotechnol.* **36**, 962 (2018).
- 6 3. Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell. *Nat. Methods* **6**, 377–382  
7 (2009).
- 8 4. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell  
9 gene expression data. *Nat. Biotechnol.* **33**, 495–502 (2015).
- 10 5. Karaiskos, N. *et al.* The Drosophila embryo at single-cell transcriptome resolution. *Science (80-. ).*  
11 **358**, 194–199 (2017).
- 12 6. Achim, K. *et al.* High-throughput spatial mapping of single-cell RNA-seq data to tissue of origin.  
13 *Nat. Biotechnol.* **33**, 503–509 (2015).
- 14 7. Halpern, K. B. *et al.* Single-cell spatial reconstruction reveals global division of labour in the  
15 mammalian liver. *Nature* **10**, (2017).
- 16 8. Faridani, O. R. & Sandberg, R. Putting cells in their place. *Nat. Biotechnol.* **33**, 490–491 (2015).
- 17 9. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* (2019).  
18 doi:10.1016/j.cell.2019.05.031
- 19 10. Welch, J. D. *et al.* Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell  
20 Identity. *Cell* **177**, 1873–1887.e17 (2019).
- 21 11. Hardoon, D. R., Szekely, S. & Shawe-Taylor, J. Canonical Correlation Analysis: An Overview with  
22 Application to Learning Methods. *Neural Comput.* **16**, 2639–2664 (2004).
- 23 12. Yang, Z. & Michailidis, G. A non-negative matrix factorization method for detecting modules in  
24 heterogeneous omics multi-modal data. *Bioinformatics* **32**, 1–8 (2016).
- 25 13. Bishop, C. M. *Pattern Recognition and Machine Learning*. (Springer, 2006).
- 26 14. Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-  
27 sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427  
28 (2018).
- 29 15. Wang, X. *et al.* Three-dimensional intact-tissue sequencing of single-cell transcriptional states.  
30 *Science* **361**, (2018).
- 31 16. Tasic, B. *et al.* Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nat.*  
32 *Neurosci.* **19**, 335–346 (2016).
- 33 17. Arias, A. M. & Hayward, P. Filtering transcriptional noise during development: Concepts and  
34 mechanisms. *Nat. Rev. Genet.* **7**, 34–44 (2006).
- 35 18. Abdi, H. & Williams, L. J. Partial Least Squares Methods: Partial Least Squares Correlation and  
36 Partial Least Square Regression. in 549–579 (2013). doi:10.1007/978-1-62703-059-5\_23
- 37 19. Berkeley Drosophila Transcription Network Project. Available at: <http://bdtnp.lbl.gov:8080/Fly>

- 1 Net/.14.
- 2 20. Luengo Hendriks, C. L. *et al.* Three-dimensional morphology and gene expression in the Drosophila  
3 blastoderm at cellular resolution I: data acquisition pipeline. *Genome Biol.* **7**, R123 (2006).
- 4 21. Bageritz, J. *et al.* Gene expression atlas of a developing tissue by single cell expression correlation  
5 analysis. *Nat. Methods* **16**, 750–756 (2019).
- 6 22. Tabata, T., Eaton, S. & Kornberg, T. B. The Drosophila hedgehog gene is expressed specifically in  
7 posterior compartment cells and is a target of engrailed regulation. *Genes Dev.* **6**, 2635–2645 (1992).
- 8 23. Patterns of gene expression in Drosophila embryogenesis.
- 9 24. Hammonds, A. S. *et al.* Spatial expression of transcription factors in Drosophila embryonic organ  
10 development. *Genome Biol.* **14**, 1–15 (2013).
- 11 25. Tomancak, P. *et al.* Global analysis of patterns of gene expression during Drosophila embryogenesis.  
*Genome Biol.* **8**, 1–24 (2007).
- 12 26. Tomancak, P. *et al.* Systematic determination of patterns of gene expression during Drosophila  
13 embryogenesis. *Genome Biol.* **3**, 1–14 (2002).
- 14 27. Clark, E. & Akam, M. Odd-paired controls frequency doubling in Drosophila segmentation by  
15 altering the pair-rule gene regulatory network. *Elife* **5**, 1–42 (2016).
- 16 28. Clark, E. *Dynamic patterning by the Drosophila pair-rule network reconciles long-germ and short-*  
17 *germ segmentation. PLoS Biology* **15**, (2017).
- 18 29. Nitzan, M., Karaiskos, N., Friedman, N. & Rajewsky, N. Gene expression cartography. *Nature* 1–6  
19 (2019). doi:10.1038/s41586-019-1773-3
- 20 30. Stuart, T. & Satija, R. Integrative single-cell analysis. *Nat. Rev. Genet.* **20**, 257–272 (2019).
- 21 31. Rood, J. E. *et al.* Toward a Common Coordinate Framework for the Human Body. *Cell* **179**, 1455–  
22 1467 (2019).
- 23

24