

# Predicting cellular position in the *Drosophila melanogaster* embryo from Single-Cell Transcriptomics data

Jovan Tanevski<sup>1,2,+</sup>, Thin Nguyen<sup>3,+</sup>, Buu Truong<sup>4,+</sup>, Nikos Karaikos<sup>5</sup>, Mehmet Eren Ahsen<sup>6,7</sup>, Xinyu Zhang<sup>8</sup>, Chang Shu<sup>8</sup>, Ke Xu<sup>8</sup>, Xiaoyu Liang<sup>8</sup>, Ying Hu<sup>8</sup>, Hoang V.V. Pham<sup>4</sup>, Li Xiaomei<sup>4</sup>, Thuc D. Le<sup>4</sup>, Adi Tarca<sup>10</sup>, Gaurav Bhatti<sup>10</sup>, Nestoras Karathanasis<sup>11</sup>, Phillip Loher<sup>11</sup>, Zhengqing Ouyang<sup>12</sup>, Yang Chen<sup>12</sup>, Disheng Mao<sup>12</sup>, Maryam Zand<sup>13</sup>, Jianhua Ruan<sup>13</sup>, Christoph Hafemeister<sup>14</sup>, Peng Qui<sup>15</sup>, Duc Tran<sup>16</sup>, Tin Nguyen<sup>16</sup>, Attila Gabor<sup>1</sup>, DREAM SCTC Consortium, Gustavo Stolovitzky<sup>17</sup>, Nikolaus Rajewsky<sup>5,\*</sup>, Julio Saez-Rodriguez<sup>1,18,\*</sup>, and Pablo Meyer<sup>17,\*</sup>

<sup>1</sup>Institute for Computational Biomedicine, Faculty of Medicine, Heidelberg University Hospital and Heidelberg University, Heidelberg, Germany

<sup>2</sup>Department of Knowledge Technologies, Jožef Stefan Institute, Ljubljana, Slovenia

<sup>3</sup>Deakin University, Geelong, Australia

<sup>4</sup>University of South Australia, Mawson Lakes, Australia

<sup>5</sup>Berlin Institute for Medical Systems Biology, Max Delbrück Center for Molecular Medicine in the Helmholtz Association, Berlin, Germany

<sup>6</sup>Icahn School of Medicine at Mount Sinai, New York City, NY, USA

<sup>7</sup>University of Illinois, Urbana-Champaign, IL, USA

<sup>8</sup>Yale University, New Haven, CT, USA

<sup>9</sup>Center for Biomedical Informatics & Information Technology, National Cancer Institute, MD, USA

<sup>10</sup>Wayne University, Detroit, MI, USA

<sup>11</sup>Computational Medicine Center, Thomas Jefferson University, Philadelphia, PA, USA

<sup>12</sup>The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA

<sup>13</sup>University of Texas at San Antonio, TX, USA

<sup>14</sup>New York Genome Center, New York City, NY, USA

<sup>15</sup>Georgia Institute of Technology, Atlanta, GA, USA

<sup>16</sup>University of Nevada, Reno, NV, USA

<sup>17</sup>IBM T.J. Watson Research Center, Yorktown Heights, NY, USA

<sup>18</sup>Joint Research Centre for Computational Biomedicine (JRC-COMBINE), Faculty of Medicine, RWTH Aachen University, Aachen, Germany

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

<sup>\*</sup>To whom correspondence should be addressed; pmeyerr@us.ibm.com

## Abstract

Single-cell sequencing technologies are rapidly evolving, in particular, although suspension single-cell RNA sequencing has become high-throughput, it loses the spatial information encoded in the position of a cell from a tissue or organism. In order to evaluate methods that reconstruct the location of single cells in the *Drosophila* embryo using single-cell transcriptomic data and the BDNP reference atlas, we organized the DREAM Single-Cell Transcriptomics challenge. Given the public availability of the ground truth, we devised a scoring and a cross-validation scheme to evaluate the soundness and robustness of the best performing algorithms. An array of methods were used by the 34 participant teams and results show that the selection of genes was essential for accurately locating the cells in the embryo. This strategy led to the identification of a set of archetypal gene expression patterns and spatial markers as participants

were able to correctly localize rare subpopulations of cells, accurately mapping both spatially restricted and scattered groups of cells. The most selected genes had a relatively high entropy and showed high spatial clustering while developmental genes such as gap and pair-ruled genes in addition to tissue defining markers were most prominent.

## 1 Introduction

1 The recent technological advances in single-cell sequencing technologies have revolutionized  
2 the biological sciences. In particular single-cell RNA sequencing (scRNAseq) methods allow  
3 transcriptome profiling in a highly parallel manner, resulting in the quantification of thousands of  
4 genes across thousands of cells of the same tissue. However, with a few exceptions [1, 2, 3, 4]  
5 current high-throughput scRNAseq methods share the drawback of losing the information relative  
6 to the spatial arrangement of the cells on the tissue during the cell dissociation step.

7 One way of regaining spatial information computationally is to appropriately combine the  
8 single-cell RNA dataset at hand with a reference database, or atlas, containing spatial expression  
9 patterns for several genes across the tissue. This approach was pursued in a few studies [5, 6, 7, 8, 9].  
10 Achim *et al* placed 139 cells using 72 reference genes with spatial information from whole mount  
11 *in situ* hybridization (WMISH) of a marine annelid and Satija *et al* developed the *Seurat* algorithm  
12 to predict position of 851 zebrafish cells based on their scRNAseq data and spatial information  
13 from *in situ*-hybridizations of 47 genes in ZFIN collection [10]. In both cases, cell positional  
14 predictions stabilized after the inclusion of 30 reference genes. Karaïkos *et al* reconstructed  
15 the early *Drosophila* embryo at single-cell resolution and while the authors were successful in  
16 their reconstruction, their approach did not lead to a predictive algorithm and mainly centered  
17 around maximizing the correlation between scRNAseq data and the expression patterns from  
18 *in situ*-hybridizations of 84 mapped genes in The Berkeley *Drosophila* Transcription Network  
19 Project (BDNTP). In this project, *in situ* hybridization data was collected resulting in a quantitative  
20 high-resolution gene expression reference atlas [11]. Indeed, Karaïkos *et al* showed that the  
21 combinatorial expression of these 84 BDNTP markers suffice to uniquely classify almost every  
22 position of the cells within the embryo.

23 In the absence of a reference database, it is also possible to regain spatial information computationally  
24 solely from the transcriptomics data by leveraging general knowledge about statistical properties  
25 of spatially mapped genes against the statistical properties of the single-cell RNA dataset  
26 [12, 13]. Bageritz *et al.* were able to reconstruct the expression map of a *Drosophila* wing disc using  
27 scRNAseq data by correlation analysis. They exploited the coexpression of non-mapping genes to  
28 a few mapping genes with known expression patterns, to predict the spatial expression patterns of  
29 824 genes [12]. Nitzan *et al.* were able to exploit the knowledge of the distribution of distances  
30 between mapping genes in physical space to predict the possible locations of cells based on the  
31 distribution of distances between genes in the expression space. They were able to successfully  
32 reconstruct the locations of cells of the *Drosophila* and zebrafish embryos from scRNAseq data [13].  
33 Although these approaches have indicated important steps to reconstruct the position of a cell in a  
34 tissue from their RNAseq expression, a global assessment is needed to evaluate the methods used  
35 and the number and nature of the genes with spatial expression information required for correctly  
36 placing cells.

37 With this purpose and to catalyze the development of new methods to predict the location  
38 of cells from scRNAseq data we organized the DREAM Single cell transcriptomics challenge.  
39 DREAM challenges are a platform for crowdsourcing collaborative competitions[14] where a  
40 rigorous evaluation of each submitted solution allows the comparison of their performance. The  
41 quality and reproducibility of each provided solutions must also be ensured. The combination of  
42 the individual solutions, i.e., the different approaches and insights to a common problem, leads

43 to an overall wisdom-of-the-crowds (WOC) solution with generally superior performance. We  
44 set up the challenge with 3 goals in mind. First, we used the data from Karaïkos *et al* to devise  
45 a variety of algorithms to test how well they could predict the localization of the cells. Second,  
46 we evaluated how the algorithms predictions depended on the number of reference genes from  
47 BDNTP with *in situ* hybridization information included in the predictions. Third, we investigated  
48 how the biological information carried in the selected genes was implemented in the algorithms to  
49 determine embryonic patterning.

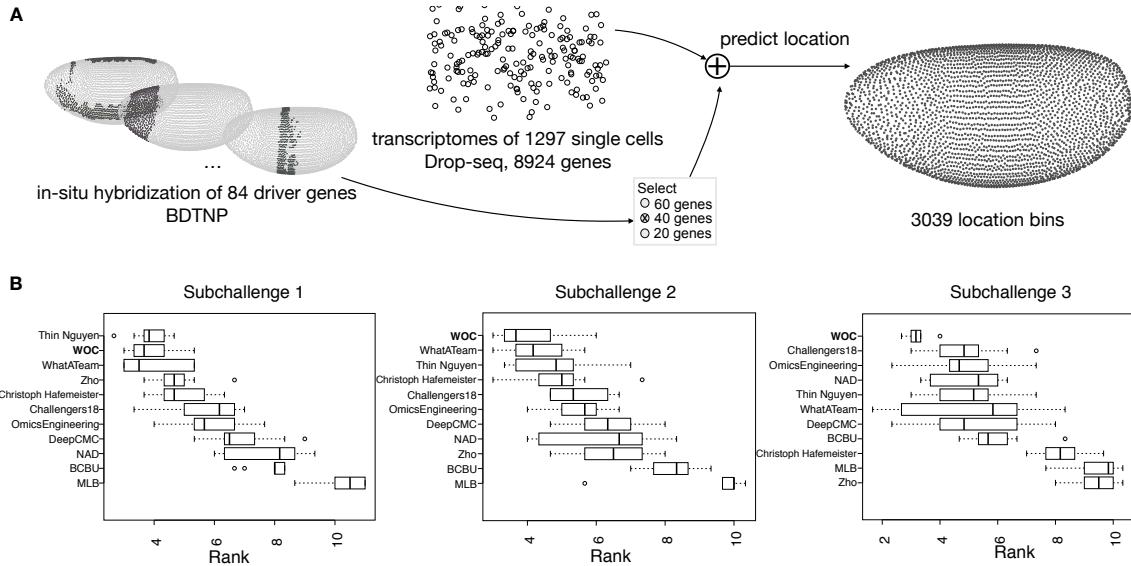
50 The challenge, a first of its kind for single cell data, consisted on predicting from their scRNAseq  
51 data the position of 1297 cells among the 3039 *Drosophila melanogaster* embryonic locations for  
52 one half of a stage 6 pre-gastrulation embryo (Figure 1A) [9]. At this stage cells in the embryo are  
53 positioned in a single two dimensional sheet following a bilateral symmetry, so that only positions  
54 in one half of the embryo were considered - accounting for the 3039 locations. Participants  
55 used the scRNAseq data for each of these 1297 cells obtained from the dissociation of 100-200  
56 stage 6 embryos and the spatial expression patterns from *in situ*-hybridizations of 84 genes in  
57 the BDNTP database [11]. Gene determinants of different tissues such as neurectoderm, dorsal  
58 ectoderm, mesoderm, yolk and pole cells were provided as a hint. To help the predictions, we  
59 provided (if they existed) the regulatory relationship -positive or negative- between the 84 genes in  
60 the *in situ*-hybridizations and the rest of the genes. We asked participants to predict the 10 most  
61 probable locations in the embryo for the 1297 cells using the expression patterns from (i) 60 genes  
62 out of the 84 in subchallenge 1, (ii) 40 genes out of the 84 in subchallenge 2, and (iii) 20 genes out  
63 of the 84 in subchallenge 3. The predictions were compared to the ground truth location using all  
64 84 *situs*. We received submissions from 34 teams, and the overall analysis of the results showed  
65 that the selection of genes is essential for accurately locating the cells in the embryo. The most  
66 selected genes had a relatively high entropy and showed high spatial clustering. In the selection,  
67 developmental genes such as gap and pair-ruled genes in addition to tissue defining markers were  
68 most prominent.

## 69 **2 Results**

### 70 **2.1 Challenge setup**

71 A distinctive feature of the single cell transcriptomics challenge was the public availability of the  
72 entire dataset and the ground truth locations produced by DistMap, a method using the *in situ*-  
73 hybridizations available at BDNTP [11], published together with the data [9]. We took three actions  
74 to mitigate the issue of the non-restricted data access. First, we prohibited the use of *in situ*s for the  
75 gene selection task and only the transcriptomics data and biological information in other databases  
76 were allowed for such task. Second, we devised three types of scores (detailed in the Methods  
77 section) that were not disclosed to the participants during the challenge. The scores measured not  
78 only the accuracy of the predicted location, but also how well the expression of the cell at the  
79 predicted location correlates to the expression from the reference atlas, the variance of the predicted  
80 locations for each cell and how well the gene-wise spatial patterns were reconstructed. Finally, we  
81 devised a post-challenge cross-validation scheme to evaluate the soundness and robustness of the  
82 methods.

83 The challenge was organized in two rounds, a leaderboard round and a final round. During the  
84 leaderboard round the participants were able to obtain scores for a limited number of solutions  
85 before submitting a single solution in the final round. We received submissions from 40 teams  
86 in the leaderboard round and 34 submissions in the final round. Out of the 34 teams that made  
87 submissions in the final round, 29 followed up with public write-ups of their approaches and



**Figure 1: Overview of the challenge and results.** **A.** In the DREAM Single-Cell Transcriptomics Challenge the participants were asked to map the location of 1297 cells to 3039 location bins of an embryo of *Drosophila melanogaster*, by combining the scRNAseq measurements of 8924 genes for each cell and the spatial expression patterns from in-situ hybridization of 60, 40 or 20 genes for each embryonic location bin, selected from a total of 84 driver-genes. **B.** Ranking of the top 10 best performing teams and a wisdom of the crowds (WOC, in bold) solution, based on results from a post challenge cross-validated selection and prediction performance measured with three complementary scoring metrics. The boxplots show the distribution of ranks for each team on the 10 test folds. The rank for each fold is calculated as the average of the ranking on each scoring metric.

source code. For subchallenges 1 and 3 we were able to determine a clear best performer, but for subchallenge 2, there were two top ranked teams with indistinguishable difference in performance (see Supplementary Figures S1,S2 and S3).

As stated, given that the dataset for this challenge was publicly available and to avoid overfitting, we decided to invite the top 10 performing teams to contribute to a post-challenge analysis phase to assess the soundness and stability of their gene selection and cell location prediction. Consequently, teams were tasked to provide predictions using a provided standard 10-fold cross-validation (CV) scenario extracted from the same dataset as in the challenge and every team used the same assignment of cells to folds. We evaluated the performance of the teams using the same scoring approach as in the challenge. To ensure the validity of the predictions, we decided to perform all the analysis and interpretation following up in this manuscript based on the 10 submitted predictions from the post-challenge phase.

## 2.2 Overview of results

Interestingly, for subchallenge 1 and 2, when participants had to use 60 or 40 genes for their predictions, the best performing teams in the CV scenario were unchanged as compared to the challenge (Figure 1B). This was not the case in subchallenge 3 as no particular team from the top 10 outperformed in a statistically significant way the others when using 20 genes for their predictions. The results from the provided cross-validation scheme showed that the approaches generalize well, i.e. the gene selection is performed consistently across the folds and the variance of the achieved scores across the folds is small for all teams (Figure S4). For each subchallenge we combined

Table 1: Best mean score for metrics  $s_1$ ,  $s_2$  and  $s_3$  achieved by the teams (Thin Nguyen, WhatATeam and OmicsEngineering) and the WOC solution. The standard deviation of scores across folds are in parenthesis. For more details on the scoring metrics see the Methods section.

	$s_1$		$s_2$		$s_3$	
	Teams	WOC	Teams	WOC	Teams	WOC
Subchallenge 1	<b>0.76(±0.04)</b>	0.72(±0.04)	<b>2.52(±0.28)</b>	2.07(±0.20)	0.59(±0.01)	<b>0.62(±0.01)</b>
Subchallenge 2	0.69(±0.03)	0.69(±0.05)	1.16(±0.12)	<b>1.82(±0.25)</b>	<b>0.67(±0.02)</b>	0.64(±0.01)
Subchallenge 3	0.65(±0.05)	<b>0.68(±0.03)</b>	0.88(±0.13)	<b>1.41(±0.16)</b>	<b>0.79(±0.02)</b>	0.70(±0.01)

108 the gene selection and location predictions from the top 10 participants into a WOC solution (see  
109 details below) that performed better compared to the individual solutions, except in subchallenge 1  
110 where two teams outperformed it (Figure 1B). The scores obtained by the best performing teams  
111 and the WOC solution are shown in Table 1.

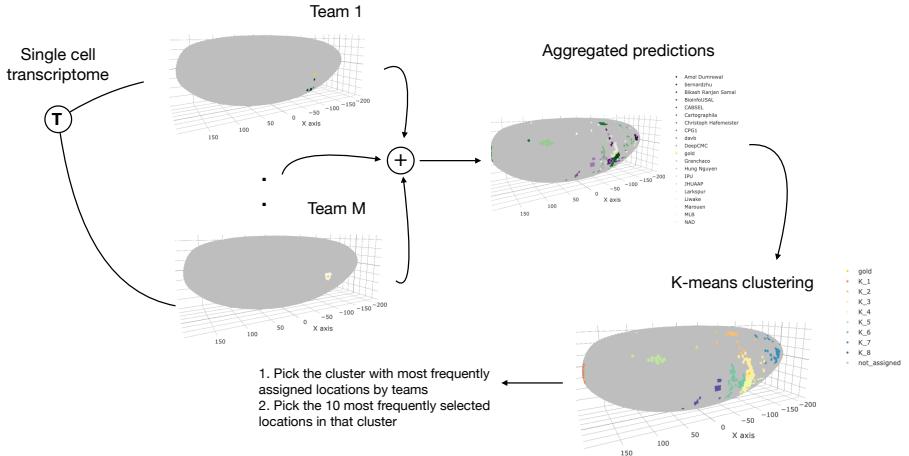
112 A summary of the methods used by participants for gene selection and location prediction can  
113 be seen in Table S1. The most frequently used method by participants for location prediction was a  
114 similarity based prediction, such as the maximum Matthews correlation coefficient between the  
115 binarized transcriptomics and the *in situ* that was proposed by Karaïkos et al. [9]. Another well  
116 performing approach was combining the predictions of a machine learning model and the Matthews  
117 correlation coefficients. The models were trained to predict either the coordinates of each cell or  
118 the binarized values of the selected *in situ* given transcriptomics data as input. The predictions  
119 were then made by selecting the location bins that corresponded to the nearest neighbors of the  
120 predicted values.

121 The most frequently used method by participants for gene selection was unsupervised or  
122 supervised feature importance estimation and ranking. For example, in a supervised feature  
123 importance estimation approach a machine learning model is trained to predict the coordinates of  
124 each cell, given the transcriptomics data at input that is the genes with available *in situ* hybridization  
125 measurements or all genes. Different machine learning models were trained by different teams.  
126 For example, Random Forest (BCBU, OmicsEngineering) or a neural network (DeepCMC, NAD).  
127 An example of an unsupervised feature importance estimation and ranking by expression based  
128 clustering (NAD, Christoph Hafemeister, MLB), or greedy feature selection based on predictability  
129 of expression from other genes (WhatATeam). Background knowledge about location specific  
130 marker genes, or the expected number of location clusters, was used by a small number of teams  
131 (WhatATeam and NAD) to inform the gene selection. Given the diversity of approaches to gene  
132 selection, we focused our analysis on better understanding the properties of frequently selected  
133 genes and providing recommendations for future experimental designs.

### 134 2.3 Analysis of the location prediction

135 A recurrent observation across DREAM challenges is that an ensemble of individual predictions  
136 performs usually better and is more robust than any individual method [15, 16]. This phenomenon,  
137 common also in other contexts, is denoted as the wisdom-of-the-crowds (WOC) [14]. In a typical  
138 challenge, individual methods output a single probability reflecting the likelihood of occurrence of  
139 an event. The WOC prediction is then constructed in an unsupervised manner by averaging the  
140 predictions of individual methods.

141 Given that in the single cell RNAseq prediction challenge participants had to submit 10 positions  
142 per cell, it did not seem consequent to use the traditional averaging of the results in the WOC  
143 approach to generate an ensemble prediction for the cell’s positions. Consequently, we developed  
144 a novel method that is based on k-means clustering for the WOC prediction. A diagram of the  
145 k-means approach is given in Figure 2. For each single cell we first used k-means clustering to



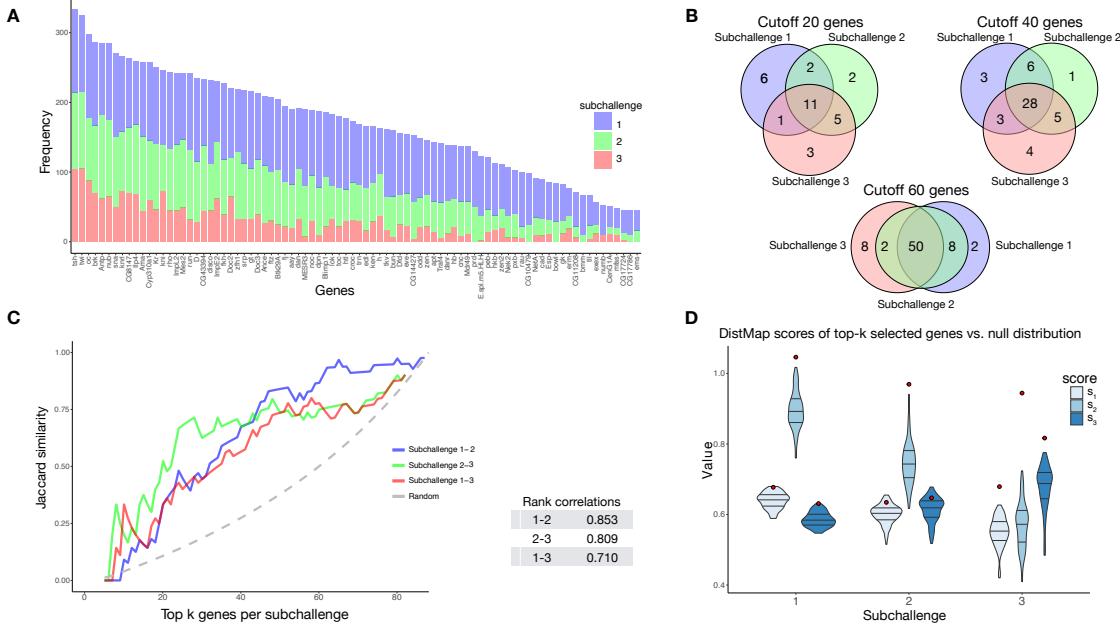
**Figure 2: Wisdom of crowds location prediction.** The location predictions for each cell by the top performing teams in the post-challenge cross-validation phase were aggregated in a wisdom of the crowds solution based on a k-means clustering approach.

cluster the locations predicted by the individual teams [17] where the euclidean distance between the locations was used as the distance for the k-means. In order to find the optimal  $k$ , we used the elbow method, i.e. we chose a  $k$  that saturates the sum of squares between clusters [18]. Note that each cluster consists of a group of locations and by construction each location is predicted by at least one team but usually more than one team. Hence, for each cluster we calculated the average frequency that its constituent locations are predicted by individual teams. We then picked the cluster with the highest average frequency as our final cluster and ranked each location in this cluster based on how frequently it was predicted by individual methods. For each site, the final prediction of the proposed WOC method consisted of the top 10 locations based on the above ranking. The k-means approach is based on the intuition that a single cell belongs to a single region and its expression is mostly similar to the locations surrounding it.

The WOC location prediction approach does not take the genes used by the teams to make the predictions into account. However, after the WOC predictions are generated, in order to score them, we need a list of genes for every subchallenge. To this end we used a WOC approach to gene selection (see the following section for more details) and used the most frequently selected genes per challenge. As reported above, the WOC solution performed better compared to the individual solutions, except in subchallenge 1 where two teams outperformed it (Figure 1B).

## 2.4 Analysis of selected genes

The selection of a subset of *in situ* used for cell location prediction was the hallmark that differentiated the subchallenges. It is unfeasible to evaluate all subsets of 20, 40 or 60 genes from the 84 due to the immense number of possible combinations of genes. Different approaches and heuristics can be used to select a subset of genes. The most frequent approaches to gene subset selection reported by the top 10 ranked teams were based on model based feature ranking algorithms, using normalized transcriptomics data (for more details see Table S1). However, if a subset of genes is selected as a candidate for solving the general task of location prediction, it should be consistent across sub-selections of cells. Therefore, we analyzed the consistency of gene selection for each team across folds by 10-fold cross-validation. More importantly, we were interested in subsets



**Figure 3: Analysis of gene selection.** The results in all figures were generated from the genes that were selected by the top performing teams in the post-challenge cross-validation scenario. **A.** Frequency of selected genes in subchallenge 1 (blue), suchallenge 2 (green) and subchallenge 3 (red). The genes are ordered according to their cumulative frequency. **B.** Venn diagrams of the most frequently selected genes in the subchallenges with cutoff at 20, 40 and 60 most frequently selected genes, corresponding to the number of genes required for each subchallenge **C.** *Left*, the similarity of most frequently selected genes for pairs of subchallenges. The Jaccard similarity measures the ratio of the size of the intersection and the union of two sets  $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ . *Right*, table of correlations between gene rankings (by frequency) for pairs of subchallenges. **D.** Validation of the performances of the wisdom of the crowds (WOC) selection of genes, i.e the most frequently selected 60, 40 and 20 genes in the respective subchallenges. The violin plots represent null distribution of scores obtained by 100 randomly selected sets of 60, 40 and 20 genes using DistMap. The red dots represents the performance obtained by using DistMap with the WOC selection of genes.

173 of genes that were consistently selected by multiple teams as this can be used as an indicator  
 174 of a biologically interesting gene or gene expression property to be used as guideline for future  
 175 experimental designs.

176 The approaches for selecting genes taken by the top 10 teams resulted in consistent selection  
 177 across folds, significantly better than random, for all subchallenges. Indeed, all of the pairwise  
 178 Jaccard similarities of sets of selected genes for all teams were significantly higher than the expected  
 179 Jaccard similarity of a random pair of subset of genes (see Supplementary Figure S4). Importantly,  
 180 we measured an observable increase in variance and decrease of mean similarity as the number of  
 181 selected genes decreased.

182 For each subchallenge we counted the number of times that the genes were selected by all teams  
 183 in all folds. The genes, ordered by the frequency of selection in all subchallenges are depicted in  
 184 Figure 3A. More than 50% of the top 20, 70% of the top 40 and 83% of the top 60 most frequently  
 185 selected genes are the same for all three subchallenges (Figure 3B). The ranks assigned to all genes  
 186 in the three subchallenges are highly correlated. Namely, the rank correlations range from 0.71  
 187 between subchallenges 1 and 3, to 0.85 between subchallenges 1 and 2. Figure 3C shows a plot of

188 the Jaccard similarity of the sets of top-k most frequently selected genes for pairs of subchallenges.  
189 We observe that a high proportion of genes are consistently selected across subchallenges. The  
190 lists of most frequently selected 60, 40 and 20 genes in subchallenges 1, 2 and 3 respectively are  
191 available in the supplementary material (Table S2).

192 We conclude that the gene selection is not only consistent by team across folds, but also across  
193 teams and subchallenges. This finding outlines a direction for further analysis, namely the validation  
194 of the predictive performance and analysis of the common properties of the most frequently selected  
195 genes.

#### 196 **2.4.1 Validation of frequently selected genes**

197 We defined a simple procedure to obtain a WOC gene selection for each of the subchallenges. It  
198 consisted on selecting the most frequently selected genes for each subchallenge (different colored  
199 bars in Figure 3A). For example, for subchallenge 1 we selected the 60 most frequently selected  
200 genes looking only at the heights of blue portion of the bar. Interestingly, the 20 most frequently  
201 selected genes in suchallenge 3 are included in the list of 40 most selected genes in subchallenge 2  
202 (except for *Doc2*), conversely included in the list of 60 most selected genes in subchallenge 1.

203 To validate the predictive performance of the WOC gene selection, we predicted the cell  
204 locations using DistMap and scored the predictions using the same scoring metrics as for the  
205 challenge, estimating the significance of the scores through generated null distributions of scores  
206 for each subchallenge. These were generated by scoring the DistMap location prediction using  
207 100 different sets of randomly selected genes. For each subchallenge and each score we estimated  
208 the empirical distribution function and then calculated the percentile of the values of the scores  
209 obtained with the WOC gene selection.

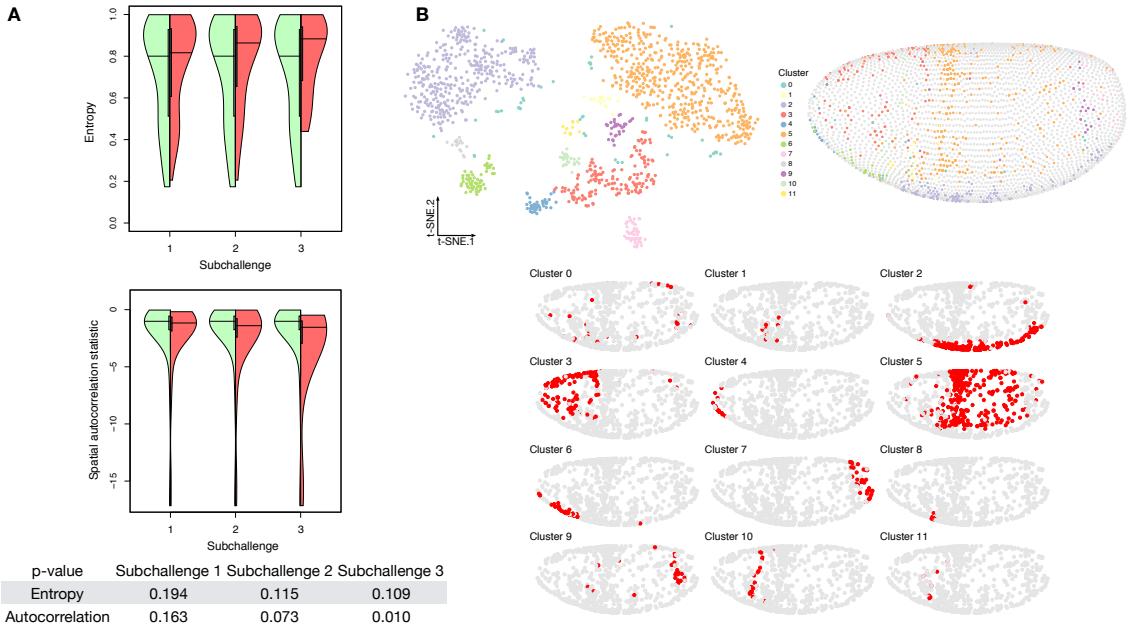
210 The null distributions and the values of the scores obtained with the WOC gene selection  
211 are shown in Figure 3D. All values of the scores fall in the 100th percentile except for  $s_1$  for  
212 subchallenge 1 (99th percentile),  $s_1$  (91st percentile) and  $s_3$  (88th percentile) for subchallenge 2.  
213 Overall the performance of DistMap with the WOC selected genes performs significantly better  
214 than a random selection of genes. The actual values of the scores are on par with those achieved by  
215 the top 10 teams in the challenge (see Table 1 for comparison).

#### 216 **2.4.2 Properties of frequently selected genes**

217 We conjectured that the most frequently selected genes should carry enough information content  
218 to uniquely encode a cell's location. Furthermore, genes should also contain location specific  
219 information, i.e. they should cluster well in space. To quantify these features, we calculated the  
220 entropy and the join count statistic for spatial autocorrelation of the *in situ* (see Figure 4A and  
221 Methods for description). We observe that most of the *in situ* genes have relatively high entropy  
222 as observed by the high density in the upper part of the plots and show high spatial clustering, i.e.  
223 show values of the join count test statistic lower than zero.

224 The Shapiro-Wilk test of normality rejected the null-hypothesis for both entropy and join count  
225 metrics ( $p < 2.3 \cdot 10^{-6}$  and  $p < 1.8 \cdot 10^{-15}$ ) for the *in situ* genes. Therefore, to test our conjectures  
226 of high entropy and spatial correlation we tested the significance of the shift of the values between  
227 all *in situ* and the WOC selected genes for each subchallenge using a one sided Mann-whitney U  
228 test. The distribution of the values for both metrics reflect that the smaller the number of selected  
229 genes, i.e going from subchallenge 1 to 3, the more significant the value shift (see Figure 4A right  
230 red part of violin plots and table).

231 To test whether the information relative to different cell types is retained with the selected  
232 subset of 60, 40 or 20 WOC selected genes, we embedded the cells into 2D space using t-distributed  
233 stochastic embedding (tSNE) [19] aiming for high accuracy ( $\theta = 0.01$ ), Figure 4B and Figure S5.



**Figure 4: Properties of selected genes.** **A.** Double violin plots of the distribution of entropy and spatial autocorrelation statistic of *Left, green* all in-situs calculated on all embryonic location bins and *Right, red* the most frequently selected 60, 40 and 20 genes in the respective subchallenges. [bottom table] p-values of a one sided Mann-whitney U test of location shift comparing the distributions on the two sides of the violin plots. **B.** *Top left*, visualization of the transcriptomics data containing only the most frequently selected 60 genes from subchallenge 1 by the top performing teams (embedding to 2D by tSNE). Each point (cell) is filled with the color of the cluster that it belongs to (density-based clustering with DBSCAN). *Top right*, spatial mapping of the cells in the *Drosophila* embryo as assigned by DistMap using only the 60 most frequently selected genes from subchallenge 1. The color of each point corresponds to the color of the cluster from the tSNE visualization. *Bottom*, highlighted (red) location mapping of cells in the *Drosophila* embryo for each cluster separately.

234 We then clustered the tSNE embedded data using density-based spatial clustering of applications  
 235 with noise (DBSCAN) [20]. DBSCAN determines the number of clusters in the data automatically  
 236 based on the density of points in space. The minimum number of cells in a local neighborhood was  
 237 set to 10 and the parameter  $\epsilon = 3.5$  was selected by determining the elbow point in a plot of sorted  
 238 distances of each cell to its 10th nearest neighbor. We found that the 9 prominent cell clusters  
 239 identified in the study by Karaïkos *et al.* [9] are preserved in our tSNE embedding and clustering  
 240 experiments when considering the most frequently selected 60 or 40 genes from subchallenges 1  
 241 and 2. The number of clusters of cells with specific localization is reduced when considering the  
 242 most frequently selected 20 genes from subchallenge 3.

243 We next associated the properties of the *in situ* that were found to be indicative of good perfor-  
 244 mance in the task of location prediction with statistical properties of the genes in the transcriptomics  
 245 data. Our goal was to discover statistical properties of the transcriptomics data that might inform  
 246 future experimental designs when selecting target genes for *in situ*-hybridizations. We calculated  
 247 statistical features across cells for the subset of genes from the transcriptomics data for which we  
 248 also have *in situ* measurements. These include the variance of gene expression  $\sigma^2$  across cells,  
 249 the coefficient of variation  $c_v = \frac{\sigma}{\mu}$ , the number of cells with expression zero 0 and the entropy of  
 250 binarized expression  $H_b$ . We then calculated the correlation across genes for each of these metrics

Table 2: Correlations of transcriptomics to *in situ* properties across most selected genes where both measurements are available.  $\sigma^2$  - variance of a gene across cells,  $c_v$  - coefficient of variation, 0 - number of cells with zero expression,  $H_b$  - entropy of binarized expression,  $H$  - entropy,  $Z$  - join count test statistic

	$\rho$	in situ	
		$H$	$Z$
scRNAseq	$\sigma^2$	0.50	0.18
	$c_v$	-0.69	0.26
	0	-0.64	0.29
	$H_b$	0.72	-0.30

and the measured spatial properties of interest of the *in situ*, i.e entropy  $H$  and the value of the joint count statistic  $Z$  (see Table 2). Although the selection of highly variable genes was one of the approaches used by some of the top 10 teams, the variance for each gene in the scRNAseq expression is surprisingly only somewhat correlated to the entropy of the corresponding *in situ* measurements of that gene. Also, we observed that the positive correlation of the entropy to the variance of each gene, becomes a negative correlation against their coefficient of variation. This negative correlation can have two sources, the genes with high entropy may have low standard deviation or high mean expression. Since we observe positive correlation of entropy to the variance of expression, we can conclude that the negative correlation is a result of highly expressed genes. Since a known drawback of scRNAseq is a high number of dropout events for lowly expressed genes [21], this observation is further supported by the negative correlation of the entropy and the number of cells with zero expression. We observed the highest correlation of *in situ* entropy to the entropy of the binarized expression.

Regarding the spatial autocorrelation, all statistical features of the transcriptomics were only slightly positively correlated to the join count statistic except for the entropy of binarized expression which had slightly negative correlation. Therefore the identification of spatially autocorrelated genes might require the use of alternative scRNAseq focused approaches [12, 13].

### 3 Discussion

In this paper we report the result of crowdsourcing as a DREAM challenge the prediction of the location of cells from scRNAseq data. Analysis of the resulting methods and their performance provided us a number of unbiased insights. First, it unveiled a connection in the cell-to-cell variability in *Drosophila* embryo gene expression and the selection of the best genes for predicting the localization of a cell in the embryo from their scRNAseq expression. The most selected genes had a relatively high entropy, hence high variance and expression while also showing high spatial clustering. The smaller the number of selected genes, i.e going from subchallenge 1 to 3, the more these features became apparent. The observed advantage of genes with high overall expression in cells might lead to less dropout counts in the scRNAseq data, a known disadvantage of the technology, leading to more accuracy in the cell placement. We also found that the 9 prominent spatially distinct cell clusters previously identified [9] are preserved when considering the most frequently selected 60 or 40 genes, but the number of clusters is reduced when considering only the most frequently selected 20 genes. This finding is in line with the conclusions of Howe et al. [10] where in a related task of location prediction the performance stabilized after the inclusion of 30 genes in a related experiment. The WOC gene selection and the k-means clustered WOC model for cell localization performed comparably or better than the participant's models, showing once more the advantage of the wisdom-of-the-crowds.

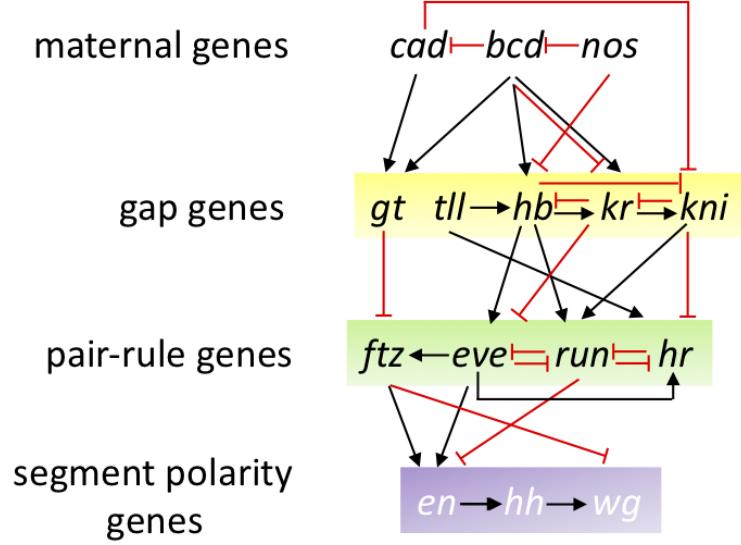


Figure 5: Gene regulatory network of early *Drosophila* development. Not all regulations are represented

Given that it has been shown that positional information of the anterior-posterior (A-P) axis is encoded as early in the embryonic development as when the expression of the gap genes occurs [22, 23], we thought that it should be possible to implement in algorithms for this challenge the information contained in the regulatory networks of *Drosophila* development [24]. Although only a small number of participants, among them the best performers, directly used biological information related to the regulation of the genes or their connectivity, the most frequently selected genes in all 3 subchallenges have interesting biological properties. Indeed, gap genes such as Giant (GT), Kruppel (KR), Knirps (KNI) were selected in all 3 subchallenges (see Fig 5 and Table S2 that also includes KNI-like KNRL) although Tailless (TLL0 and Hunchback (HB) were not. Along the A-P axis, maternally provided Bicoid (BCD) and Caudal (CAD) first establish the expression patterns of gap and terminal class factors, such as HB, GT, KR and KNI. These A-P early regulators then collectively direct transcription of A-P pair-rule factors, such as even-skipped (EVE), Fushi-tarazu (FTZ), Hairy (H), Odd skipped, (ODD), Paired (PRD) and Runt (RUN) which in turn cross-regulate each other. Not being part of the *in situ*, neither BCD, nor CAD were selected but AMA sitting near BCD in the genome might have been selected for its similar expression properties. Furthermore, we also find that pair-rule genes were most prominently selected in subchallenges 1 ( EVE, ODD, PRD, the Paired-like NOC and RUN) and 2 (H, FTZ and RUN). A similar cascade of maternal and zygotic factors controls patterning along the dorsal-ventral (D-V) axis where Dorsal (D), Snail (SNA) and Twist (TWI) specify mesoderm and the pair rule factors EVE and FTZ specify location along the trunk of the A-P axis. Again, SNA and TWI were selected in all subchallenges and D in subchallenges 1 and 2. These selected transcription factors specify distinct developmental fates and can act via different cis-regulatory modules but their quantitative differences in relative levels of binding to shared targets correlates with their known biological and transcriptional regulatory specificities [25]. The rest of the selected genes were the homeobox genes (NUB, ANTP) and differentiators of tissue such as mesoderm (AMA, MES2, ZFH1), ectoderm (DOC2 and DOC3), neural tissue (NOC, OC, RHO) and EGFR pathway (RHO, EDL). The complete lists of most frequently selected genes are available in Table S2.

313 Since the entire dataset and the ground truth locations were publicly available, the organization  
314 of this DREAM challenge brought risks that, given the importance of the scientific question asked,  
315 we thought worth taking. However, without the post-challenge phase it would have been impossible  
316 to distinguish the robust and sound methods from methods that were overfitting the results. Overall,  
317 the single cell transcriptomics challenges unveils not only the best gene-selection methods and  
318 prediction approaches to place a cell in the *Drosophila* embryo, but also explains the biological and  
319 statistical properties of the genes selected for the predictions. We think that such properties could  
320 be used or adapted when performing similar cell-placing tasks in other organisms, including human  
321 tissues. Given the importance of spatial arrangements for disease development and treatment, we  
322 foresee an application of these methods to medical questions as well.

## 323 References

- 324 [1] Anna K Casasent, Aislyn Schalck, Ruli Gao, Emi Sei, Annalyssa Long, William Pangburn,  
325 Tod Casasent, Funda Meric-Bernstam, Mary E Edgerton, and Nicholas E Navin. Multiclonal  
326 invasion in breast tumors identified by topographic single cell sequencing. *Cell*, 172(1-2):205–  
327 217, 2018.
- 328 [2] Patrik L Ståhl, Fredrik Salmén, Sanja Vickovic, Anna Lundmark, José Fernández Navarro,  
329 Jens Magnusson, Stefania Giacomello, Michaela Asp, Jakub O Westholm, Mikael Huss, et al.  
330 Visualization and analysis of gene expression in tissue sections by spatial transcriptomics.  
331 *Science*, 353(6294):78–82, 2016.
- 332 [3] Ditte Lovatt, Brittani K Ruble, Jaehee Lee, Hannah Dueck, Tae Kyung Kim, Stephen Fisher,  
333 Chantal Francis, Jennifer M Spaethling, John A Wolf, M Sean Grady, et al. Transcriptome in  
334 vivo analysis (tiva) of spatially defined single cells in live tissue. *Nature methods*, 11(2):190,  
335 2014.
- 336 [4] Samuel G Rodrigues, Robert R Stickels, Aleksandrina Goeva, Carly A Martin, Evan Murray,  
337 Charles R Vanderburg, Joshua Welch, Linlin M Chen, Fei Chen, and Evan Z Macosko. Slide-  
338 seq: A scalable technology for measuring genome-wide expression at high spatial resolution.  
339 *Science*, 363(6434):1463–1467, 2019.
- 340 [5] Kaia Achim, Jean-Baptiste Pettit, Luis R Saraiva, Daria Gavriouchkina, Tomas Larsson,  
341 Detlev Arendt, and John C Marioni. High-throughput spatial mapping of single-cell rna-seq  
342 data to tissue of origin. *Nature biotechnology*, 33(5):503, 2015.
- 343 [6] Rahul Satija, Jeffrey A Farrell, David Gennert, Alexander F Schier, and Aviv Regev. Spatial  
344 reconstruction of single-cell gene expression data. *Nature biotechnology*, 33(5):495, 2015.
- 345 [7] Keren Bahar Halpern, Rom Shenhav, Orit Matcovitch-Natan, Beáta Tóth, Doron Lemze,  
346 Matan Golan, Efi E Massasa, Shaked Baydatch, Shanie Landen, Andreas E Moor, et al.  
347 Single-cell spatial reconstruction reveals global division of labour in the mammalian liver.  
348 *Nature*, 542(7641):352, 2017.
- 349 [8] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi,  
350 William M Mauck III, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija.  
351 Comprehensive integration of single-cell data. *Cell*, 2019.
- 352 [9] Nikos Karaiskos, Philipp Wahle, Jonathan Alles, Anastasiya Boltengagen, Salah Ayoub,  
353 Claudia Kipar, Christine Kocks, Nikolaus Rajewsky, and Robert P Zinzen. The drosophila  
354 embryo at single-cell transcriptome resolution. *Science*, 358(6360):194–199, 2017.

- 355 [10] Douglas G Howe, Yvonne M Bradford, Tom Conlin, Anne E Eagle, David Fashena, Ken  
356 Frazer, Jonathan Knight, Prita Mani, Ryan Martin, Sierra A Taylor Moxon, et al. Zfin, the  
357 zebrafish model organism database: increased support for mutants and transgenics. *Nucleic  
358 acids research*, 41(D1):D854–D860, 2012.
- 359 [11] Charless C Fowlkes, Cris L Luengo Hendriks, Soile VE Keränen, Gunther H Weber, Oliver  
360 Rübel, Min-Yu Huang, Sohail Chattoor, Angela H DePace, Lisa Simirenko, Clara Henriquez,  
361 et al. A quantitative spatiotemporal atlas of gene expression in the drosophila blastoderm.  
362 *Cell*, 133(2):364–374, 2008.
- 363 [12] Josephine Bageritz, Philipp Willnow, Erica Valentini, Svenja Leible, Michael Boutros, and  
364 Aurelio A Teleman. Gene expression atlas of a developing tissue by single cell expression  
365 correlation analysis. *Nature Methods*, 16(8):750, 2019.
- 366 [13] Mor Nitzan, Nikos Karaikos, Nir Friedman, and Nikolaus Rajewsky. Charting a tissue from  
367 single-cell transcriptomes. *bioRxiv*, page 456350, 2018.
- 368 [14] Julio Saez-Rodriguez, James C Costello, Stephen H Friend, Michael R Kellen, Lara Man-  
369 gravite, Pablo Meyer, Thea Norman, and Gustavo Stolovitzky. Crowdsourcing biomedical  
370 research: leveraging communities as innovation engines. *Nature Reviews Genetics*, 17(8):470–  
371 486, 2016.
- 372 [15] Daniel Marbach, James C Costello, Robert Küffner, Nicole M Vega, Robert J Prill, Diogo M  
373 Camacho, Kyle R Allison, Andrej Aderhold, Richard Bonneau, Yukun Chen, et al. Wisdom  
374 of crowds for robust gene network inference. *Nature methods*, 9(8):796, 2012.
- 375 [16] Michael P Menden, Dennis Wang, Mike J Mason, Bence Szalai, Krishna C Bulusu, Yuanfang  
376 Guan, Thomas Yu, Jaewoo Kang, Minji Jeon, Russ Wolfinger, et al. Community assessment  
377 to advance computational prediction of cancer drug combinations in a pharmacogenomic  
378 screen. *Nature communications*, 10(1):2674, 2019.
- 379 [17] John A Hartigan and Manchek A Wong. Algorithm as 136: A k-means clustering algorithm.  
380 *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):100–108, 1979.
- 381 [18] Trupti M Kodinariya and Prashant R Makwana. Review on determining number of cluster in  
382 k-means clustering. *International Journal*, 1(6):90–95, 2013.
- 383 [19] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of  
384 Machine Learning Research*, 9:2579–2605, 2008.
- 385 [20] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm  
386 for discovering clusters in large spatial databases with noise. In *Proceedings of the Second  
387 International Conference on Knowledge Discovery and Data Mining*, KDD’96, pages 226–  
388 231, 1996.
- 389 [21] Peter V Kharchenko, Lev Silberstein, and David T Scadden. Bayesian approach to single-cell  
390 differential expression analysis. *Nature methods*, 11(7):740, 2014.
- 391 [22] Julien O Dubuis, Gašper Tkačik, Eric F Wieschaus, Thomas Gregor, and William Bialek. Posi-  
392 tional information, in bits. *Proceedings of the National Academy of Sciences*, 110(41):16301–  
393 16308, 2013.
- 394 [23] Mariela D Petkova, Gašper Tkačik, William Bialek, Eric F Wieschaus, and Thomas Gregor.  
395 Optimal decoding of cellular identities in a genetic network. *Cell*, 176(4):844–855, 2019.

- 396 [24] David Umulis, Michael B O'Connor, and Hans G Othmer. Robustness of embryonic spatial  
 397 patterning in *drosophila melanogaster*. *Current topics in developmental biology*, 81:65–111,  
 398 2008.
- 399 [25] Stewart MacArthur, Xiao-Yong Li, Jingyi Li, James B Brown, Hou Cheng Chu, Lucy Zeng,  
 400 Brandi P Grondona, Aaron Hechmer, Lisa Simirenko, Soile VE Keränen, et al. Developmental  
 401 roles of 21 *drosophila* transcription factors are determined by quantitative differences in  
 402 binding to an overlapping set of thousands of genomic regions. *Genome biology*, 10(7):R80,  
 403 2009.
- 404 [26] Andrew D. Cliff and John K. Ord. *Spatial autocorrelation*. Pion London, 1973.
- 405 [27] Robert R. Solak and Neal L. Oden. Spatial autocorrelation in biology: 1. methodology.  
 406 *Biological Journal of the Linnean Society*, 10(2):199–228, 1978.

## 4 Methods

### 4.1 Scoring

409 We scored the submissions for the three subchallenges using three metrics  $s_1$ ,  $s_2$  and  $s_3$ .  $s_1$  measured  
 410 how well the expression of the cell at the predicted location correlates to the expression from the  
 411 reference atlas and included the variance of the predicted locations for each cell. While  $s_2$  measured  
 412 the accuracy of the predicted location and  $s_3$  measured how well the gene-wise spatial patterns  
 413 were reconstructed.

414 Let  $c$  represent the identity of a cell, given in the transcriptomics data in the challenge where  
 415  $1 \leq c \leq 1297$ . Each cell  $c$  is located in a bin  $\varepsilon_c \in \{1..3039\}$  at a position with coordinates  
 416  $r(\varepsilon_c) = (x_c, y_c, z_c)$ . Each cell is associated with a binarized expression profile  $t_c = (t_{c1}, t_{c2}, \dots, t_{cE})$ ,  
 417 where  $1 \leq E \leq 3100$ , and a corresponding binarized *in situ* profile  $f_c = (f_{c1}, f_{c2}, \dots, f_{cK})$ , where  
 418 the maximum possible value of  $K$  for which we have *in situ* information is  $K = 84$ . For different  
 419 subchallenges we consider  $K \in \{20, 40, 60\}$ . Using  $K$  selected genes the participants were asked to  
 420 provide an ordered list of 10 most probable locations for each cell. We represent with the mapping  
 421 function  $A(c, i, K)$  the value of the predicted  $i$ -th most probable location for cell  $c$  using  $K$  *in situ*.

422 For the first scoring metric  $s_1$  we calculated the weighted average of the Mathews correlation  
 423 coefficient (MCC) between the *in situ* profile of the ground truth cell location  $f_{\varepsilon_c}$  and the *in situ*  
 424 profile of the most probable predicted location  $f_{A(c, 1, K)}$  for that cell

$$s_1 = \sum_{c=1}^N \frac{p_K(c, A)}{\sum_{c=1}^N p_K(c, A)} MCC(f_{A(c, 1, K)}, f_{\varepsilon_c}).$$

425 The Matthews correlation coefficient, or  $\phi$  coefficient, is calculated from the contingency table  
 426 obtained by comparing two binary vectors as  $MCC = \frac{ad - bc}{\sqrt{(a+b)(a+c)(d+b)(d+c)}}$ , where  $a$  is the number  
 427 of true positives,  $b$  is the number of false positives,  $c$  the number of false negatives and  $d$  the number  
 428 of true negatives. The MCC is weighted by the inverse of the distance of the predicted most probable  
 429 locations to the ground truth location  $p_K(c)$ . The weights are calculated as  $p_K(c, A) = \frac{d_{84}(c, A)}{d_K(c, A)}$ ,  
 430 where  $d_K(c, A) = \frac{1}{10} \sum_{i=1}^{10} \|r(A(c, i, K)) - r(\varepsilon_c)\|_2$ ,  $d_{84}(c, A)$  is the value of  $d_K(c, A)$  using the ground  
 431 truth most probable locations assigned with  $K = 84$  using DistMap, and  $\|\cdot\|_2$  is the Euclidean norm.

The second metric  $s_2$  is simply the average inverse distance of the predicted most probable  
 locations to the ground truth location

$$s_2 = \frac{1}{N} \sum_{c=1}^N p_K(c, A).$$

432 Finally, the third metric  $s_3$  measures the accuracy of reconstructed gene-wise spatial patterns

$$s_3 = \sum_{s=1}^K \frac{MCC(t_{cs}, f_{\epsilon_{cs}})_{\forall c}}{\sum_{s=1}^K MCC(t_{cs}, f_{\epsilon_{cs}})_{\forall c}} MCC(t_{cs}, f_{A(c,1,K)s})_{\forall c}.$$

433 For 287 out of the 1297 cells, the ground truth location predictions were ambiguous, i.e., the  
 434 MCC scores were identical for multiple locations. These cells were removed both from the ground  
 435 truth and the submissions before calculating the scores.

436 The teams were ranked according to each score independently. The final assigned rank  $r_t$   
 437 for team  $t$  was calculated as the average rank across scores. Teams were ranked based on the  
 438 performance as measured by the three scores on 1000 bootstrap replicates of the submitted solutions.  
 439 The three scores were calculated for each bootstrap. The teams were then ranked according to  
 440 each score. These ranks were then averaged to obtain a final rank for each team on that bootstrap.  
 441 The winner for each subchallenge was the team that achieved the lowest ranks. We calculated the  
 442 Bayes factor of the bootstrap ranks for the top performing teams. Bayesian factor of 3 or more was  
 443 considered as a significantly better performance. The Bayes factor of the 1000 bootstrapped ranks  
 444 of teams  $r_1$  and  $r_2$  was calculated as

$$BF(r_1, r_2) = \frac{\sum_{i=1}^{1000} \mathbf{1}(r_{1\_i} < r_{2\_i})}{\sum_{i=1}^{1000} \mathbf{1}(r_{1\_i} > r_{2\_i})},$$

445 where  $\mathbf{1}$  is the indicator function.

## 446 4.2 Entropy and spatial autocorrelation

The entropy of a binarized *in situ* measurements of gene  $G$  was calculated as

$$H(G) = -p \log_2 p - (1-p) \log_2 (1-p),$$

447 where  $p$  is the probability of gene  $G$  to have value 1.

448 The join count statistic is a measure of a spatial auto-correlation of a binary variable. Let  $n_B$  be  
 449 the number of bins where  $G$  is expressed ( $G = B$ ), and  $n_W = n - n_B$  the number of bins where  $G$  is  
 450 not expressed ( $G = W$ ). Two neighboring spatial bins can form join of type  $J \in \{WW, BB, BW\}$ .

451 We are interested in the distribution of BW joins. If a gene has a lower number of BW joins than  
 452 the expected number of BW, then the gene is positively spatially auto-correlated, i.e., the gene is  
 453 highly clustered. Contrarily, higher number of BW joins points towards negative spatial correlation,  
 454 i.e. dispersion.

Following Cliff and Ord [26] and Sokal and Oden [27], the expected count of BW joins is

$$\mathbb{E}[BW] = \frac{1}{2} \sum_i \sum_j \frac{w_{ij} n_B^2}{n^2},$$

where the spatial connectivity matrix  $w$  is defined as

$$w_{ij} = \begin{cases} 1 & \text{if } i \neq j \text{ and } j \text{ is in the list of 10 nearest neighbors of } i \\ 0 & \text{otherwise} \end{cases}$$

The variance of BW joins is

$$\sigma_{BW}^2 = \mathbb{E}[BW^2] - \mathbb{E}[BW]^2.$$

where the term  $\mathbb{E}[BW^2]$  is calculated as

$$\mathbb{E}[BW^2] = \frac{1}{4} \left( \frac{2x_2 n_B n_W}{n^2} + \frac{(x_3 - 2x_2)n_B n_W(n_B + n_W - 2)}{n^3} + \frac{4(x_1^2 + x_2 - x_3)n_B^2 n_W^2}{n^4} \right),$$

455 where  $x_1 = \sum_i \sum_j w_{ij}$ ,  $x_2 = \frac{1}{2} \sum_i \sum_j (w_{ij} - w_{ji})^2$ ,  $x_3 = \sum_i (\sum_j w_{ij} + \sum_j w_{ji})^2$ .

456 Note that the connectivity matrix  $w$  can also be asymmetric, since it is defined by the nearest  
457 neighbor function.

Finally, the observed BW counts are

$$BW = \frac{1}{2} \sum_i \sum_j w_{ij} (G_i - G_j)^2.$$

The join counts test statistic is then defined as

$$Z(BW) = \frac{BW - \mathbb{E}[BW]}{\sqrt{\sigma_{BW}^2}},$$

458 which is assumed to be asymptotically normally distributed under the null hypothesis of no spatial  
459 autocorrelation. Negative values of the Z statistic represent positive spatial autocorrelation, or  
460 clustering, of gene  $G$ . Positive values of the Z statistic represent negative spatial autocorrelation, or  
461 dispersion, of gene  $G$ .

### 462 4.3 Implementation details

463 The challenge scoring was implemented and run in R version 3.5, the post analysis was performed  
464 with R version 3.6 and the core `tidyverse` packages. We used the publicly available implemen-  
465 tation of `DistMap` (<https://github.com/rajewsky-lab/distmap>). MCC calculated  
466 with R package `mccr` (0.4.4). tSNE embedding and visualization produced with R package `Rtsne`  
467 (0.15). DBSCAN cluterung with R package `dbscan` (1.1-3).

### 468 4.4 Code availability

469 <https://github.com/dream-sctc/Scoring>

### 470 4.5 Data availability

#### 471 Reference Database

472 The reference database comes from the Berkeley *Drosophila* Transcription Network Project.  
473 The *in situ* expression of 84 genes (columns) is quantified across the 3039 *Drosophila* embryonic  
474 locations (rows) for raw data and for binarized data. The order of the rows in both files follows the  
475 order of the coordinates in the `geometry.txt` file. The 84 genes were binarized by manually choosing  
476 thresholds for each gene.

#### 477 Spatial coordinates

478 One half of *Drosophila* embryo has 3039 cells places as x, y and z (columns) for a total of 3039  
479 embryo locations (rows) and a total of 3039·3 coordinates in the file `geometry.txt`.

#### 480 Single cell RNA sequencing

481 The single-cell RNA sequencing data is provided as a matrix with 8924 genes as rows and  
482 1297 cells as columns. In the raw version of the matrix, the entries are the raw unique gene counts  
483 (quantified by using unique molecular identifiers – UMI). The normalized version is obtained by  
484 dividing each entry by the total number of UMIs for that cell, adding a pseudocount and taking  
485 the logarithm of that. All entries are finally multiplied by a constant. The normalized data can be  
486 obtained in `dgenormalizedtxt`, the raw data in `dgerawtxt`. For a given gene and only considering the  
487 Drop-seq cells expressing it we computed a quantile value above (below) which the gene would be  
488 designated ON (OFF). We sampled a series of quantile values and each time the gene correlation  
489 matrix based on this binarized version `dgebinarizeddistMapcsv` of our `dgenormalizedtxt` versus the

490 binarized BDTNP atlas was computed and compared by calculating the mean square root error  
491 between the elements of the lower triangular matrices. Eventually, we selected the quantile value  
492 0.23, as it was found to minimize the distance between the two correlation matrices.

493 The short sequences for each of the 1297 cells in the raw and normalized data are the cell  
494 barcodes.

## 495 **5 Acknowledgments**

496 This research was funded in part by

## 497 **6 Author contributions**

498 Conceptualization, N.K., N.R., J.S.R., G.S., and P.M.; Methodology, J.S., M.E.A., G.S., and P.M.;  
499 Software, J.T., and M.E.A.; Formal Analysis, J.T., M.E.A., G.S. and P.M.; Writing - Original Draft,  
500 J.T. and P.M.; Writing - Supervision, J.S.R., G.S., and P.M.

## 501 **7 Competing interests**

502 The authors declare no competing interests.

## 503 **8 Materials and Correspondence**

504 Requests for data, resources, and or reagents should be directed to Pablo Meyer (pmeyerr@us.ibm.com).

505 **Supplementary information**

Subchallenge 1

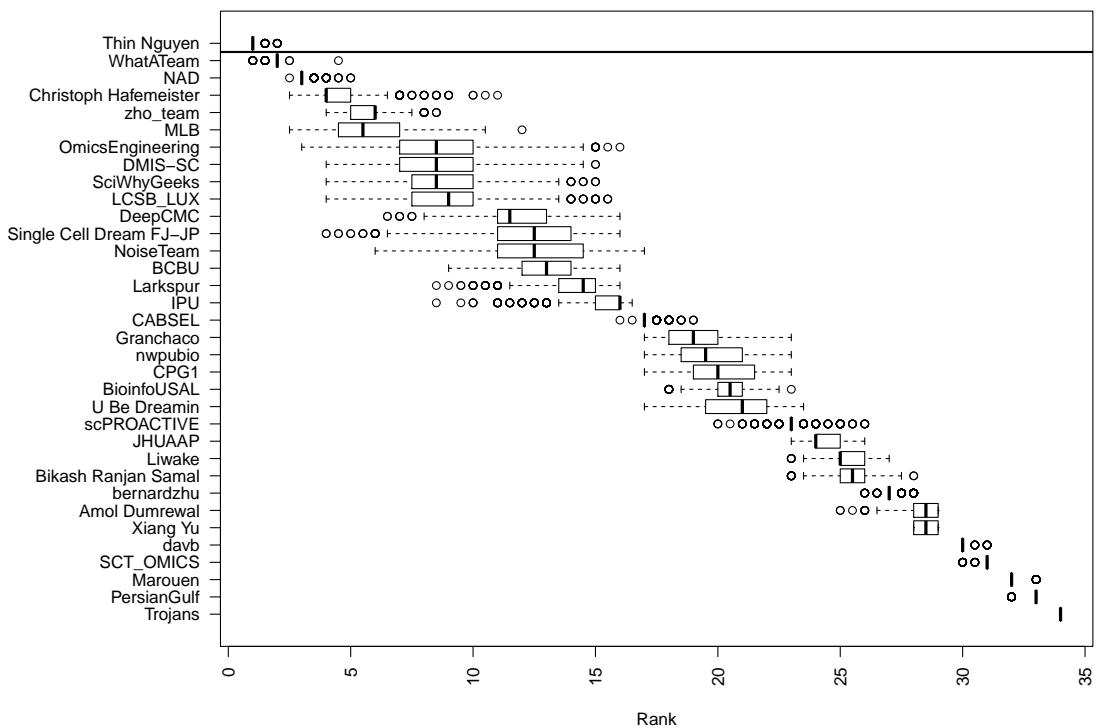


Figure S1: Results from the challenge

## Subchallenge 2

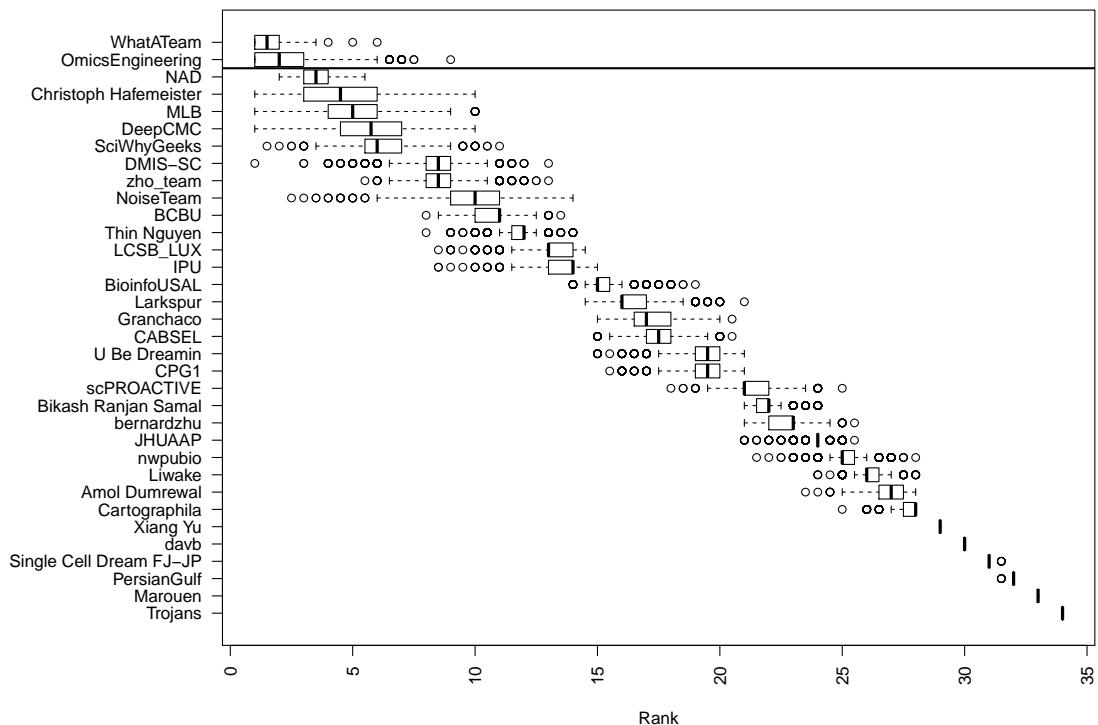


Figure S2: Results from the challenge

### Subchallenge 3

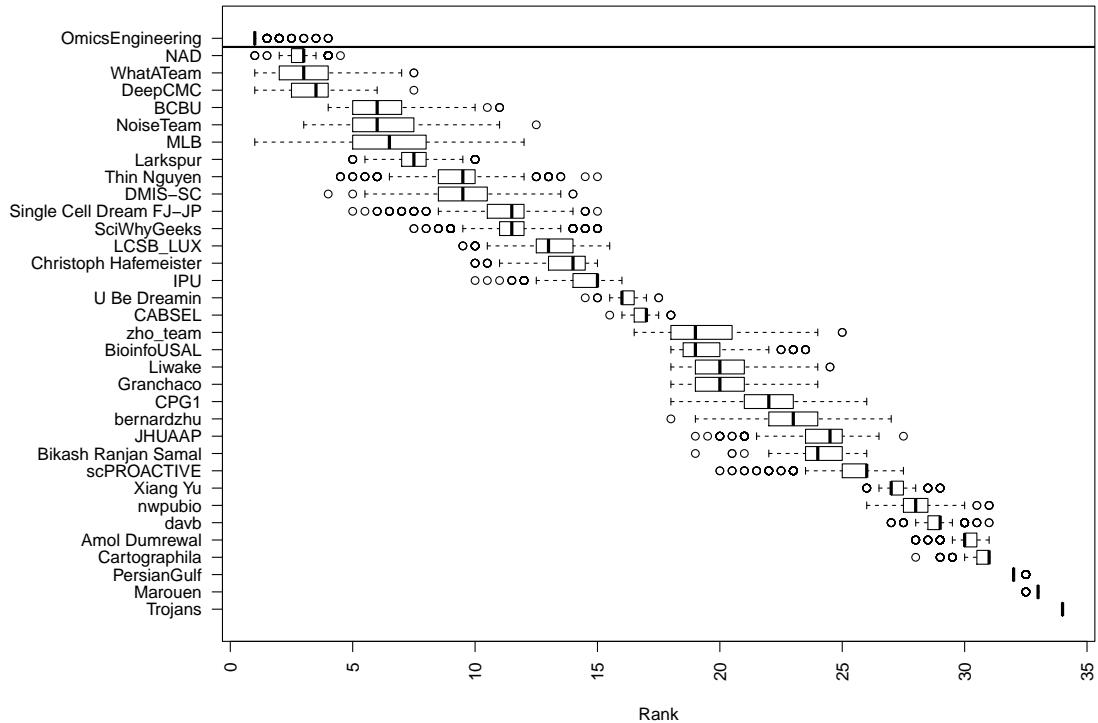


Figure S3: Results from the challenge

Table S1: Methods used by the top 10 teams for gene selection and location prediction. Some teams used different approaches or a combination of approaches for different subchallenges. SFR - Supervised feature ranking, UFR - unsupervised feature ranking, KNW - background knowledge, VAR - variance. CMB - Combination of model prediction and MCC, MCC - Matthews correlation coefficient, SIM - Similarity measure (non MCC)

Prediction	Selection			
	SFR	UFR	KNW	VAR
CMB	2	1	2	
MCC	3	2		1
SIM			5	

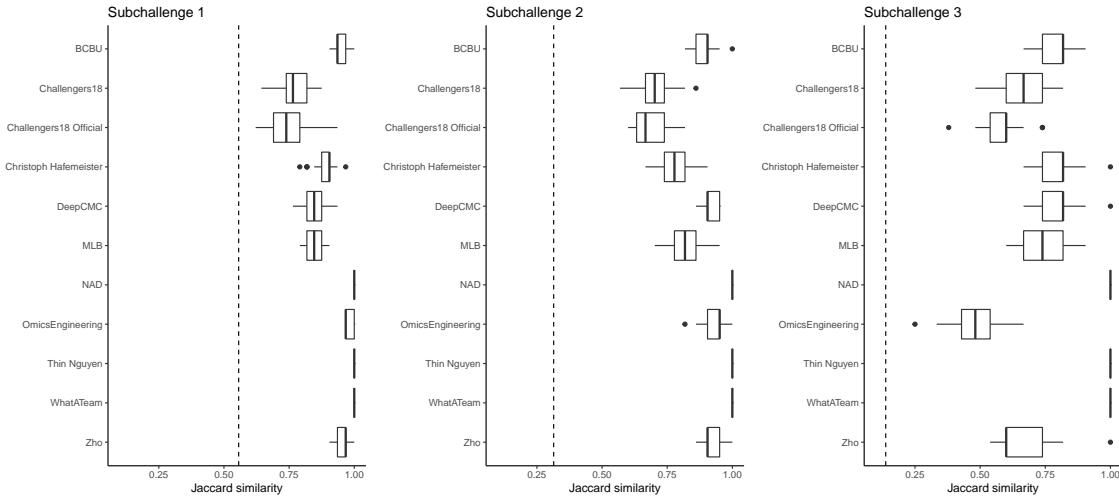


Figure S4: Boxplots of the Jaccard similarity between the genes selected for each of the 10 CV scheme in all 3 subchallenges. The teams that used the statistical properties of the genes as selection criteria, for example maximum variance, selected the same set of genes for all folds. This is expected since the distribution of a random subsample was selected to have the same properties as the original sample. Dotted line represents the limit for significance, i.e., the expected Jaccard similarity between two sets of randomly selected 60, 40 or 20 genes.

Table S2: Most frequently selected 60, 40 and 20 genes in subchallenges 1,2 and 3 respectively, in alphabetical order, colored according to Figure 5 from the main text.

Subchallenge 1	aay Ama Ance Antp apt Blimp-1 brk Btk29A bun CG14427 CG43394 CG8147 cnc croc Cyp310a1 D dan danr Dfd disco Doc2 Doc3 dpn E(spl)m5-HLH edl eve fj fkh ftz gt htl Ilp4 ImpE2 ImpL2 ken kni knrl Kr lok Mdr49 Mes2 MESR3 noc nub oc odd prd rau rho run sna srp tkv toc Traf4 trn tsh twi zen zfh1
Subchallenge 2	aay Ama Ance Antp Blimp-1 brk Btk29A bun CG43394 CG8147 Cyp310a1 D dan disco Doc3 dpn edl fj fkh ftz gt h Ilp4 ImpE2 ImpL2 kni knrl Kr Mes2 MESR3 noc nub oc rho run sna srp tsh twi zfh1
Subchallenge 3	Ama Antp brk CG8147 Cyp310a1 disco Doc2 Ilp4 ImpE2 ImpL2 kni knrl Kr Mes2 nub oc rho sna tsh twi

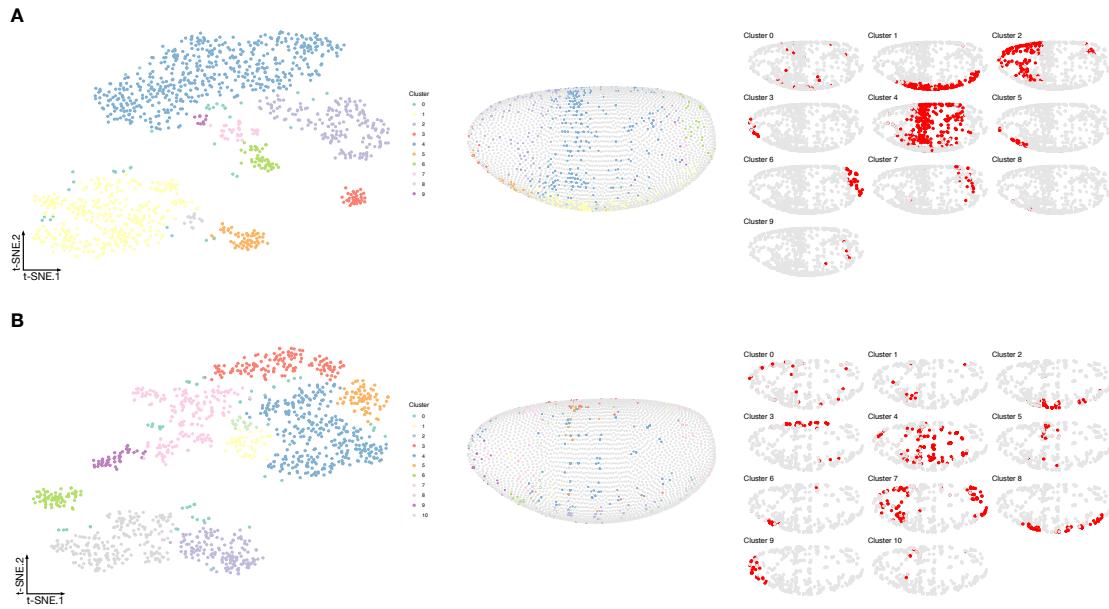


Figure S5: Visualization of the transcriptomics data containing only the most frequently selected **A** 40 genes from subchallenge 2 and **B** 20 genes from subchallenge 3 by the top performing teams (embedding to 2D by tSNE). *Left* each point (cell) is filled with the color of the cluster that it belongs to (density-based clustering with DBSCAN). *Middle*, spatial mapping of the cells in the Drosophila embryo as assigned by DistMap using only the 60 most frequently selected genes from subchallenge 1. The color of each point corresponds to the color of the cluster from the tSNE visualization. *Right*, highlighted (red) location mapping of cells in the Drosophila embryo for each cluster separately.