

# <sup>1</sup> ELeFHAnt: Ensemble Learning for Harmonization and Annotation of Single Cell Data

<sup>3</sup> Konrad Thorner<sup>1</sup> and Praneet Chaturvedi<sup>1</sup>

<sup>4</sup> 1 Cincinnati Children's Hospital Medical Center

DOI: [10.21105/joss.03516](https://doi.org/10.21105/joss.03516)

## Software

- [Review ↗](#)
- [Repository ↗](#)
- [Archive ↗](#)

---

Editor: Pending Editor ↗

Submitted: 19 July 2021

Published: 21 July 2021

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

## <sup>5</sup> Summary

<sup>6</sup> Single cell sequencing has become a powerful method for understanding biological systems at an increasingly granular level. For single cell RNA-seq data specifically, one of the primary <sup>7</sup> questions is using the transcriptome to determine cell identity. It is common to visualize such <sup>8</sup> data and find clusters of cells with similarity in gene expression, but assigning each cluster a <sup>9</sup> cell type is a much more open-ended task. Taking advantage of publicly-available, annotated <sup>10</sup> datasets in combination with supervised learning is a powerful method for addressing this <sup>11</sup> question. Ensemble Learning for Harmonization and Annotation of Single Cells (ELeFHAnt) <sup>12</sup> provides an easy to use R package for users to annotate clusters of single cells, harmonize <sup>13</sup> labels across single cell datasets to generate a unified atlas, and infer relationships among cell <sup>14</sup> types between two datasets. It provides users with the flexibility of choosing between random <sup>15</sup> forest and SVM (Support Vector Machine) based classifiers or letting ELeFHAnt apply both <sup>16</sup> in combination to make predictions.

## <sup>18</sup> Statement of Need

<sup>19</sup> As an alternative to manual annotation, there are many automatic cell annotation tools currently available that employ either gene marker, correlation, or machine learning-based methods, each with varying levels of performance ([Pasquini et al., 2021](#)). The label transfer <sup>20</sup> functionality of Seurat is among the most well-known, but occurs on an individual cell level <sup>21</sup> rather than a community level, often leading to over-annotation ([Stuart et al., 2019](#)). There <sup>22</sup> are also deep learning-based tools emerging such as scANVI, which utilizes generative models <sup>23</sup> but requires significantly more computation time ([Xu et al., 2021](#)).

<sup>26</sup> ELeFHAnt is a supervised machine learning-based tool that enables researchers to identify cell <sup>27</sup> types in their scRNA-seq data while providing additional unique features. ELeFHAnt gives <sup>28</sup> users the ability to use and compare not just one but multiple classification algorithms simultaneously through its ensemble method, weighting their predictions to produce the best <sup>29</sup> consensus among them. In this ensemble, SVM and random forest are our two classifiers selected <sup>30</sup> for their superior accuracy and computation time in a benchmarking study ([Abdelaal et al., 2019](#)). Additionally, selecting the optimal reference is a challenge addressed by harmonization, <sup>31</sup> that allows users to integrate multiple datasets together into an atlas. A standardized <sup>32</sup> set of labels is generated across all of them, which can subsequently be used to annotate new <sup>33</sup> datasets. Relationships between two datasets can also be deduced to better understand how <sup>34</sup> each was annotated. This is provided in an easy to interpret heatmap format that compares <sup>35</sup> all the cell types between them. Finally, a subsampling procedure is used to enable faster <sup>36</sup> predictions while being shown not to influence reproducibility.

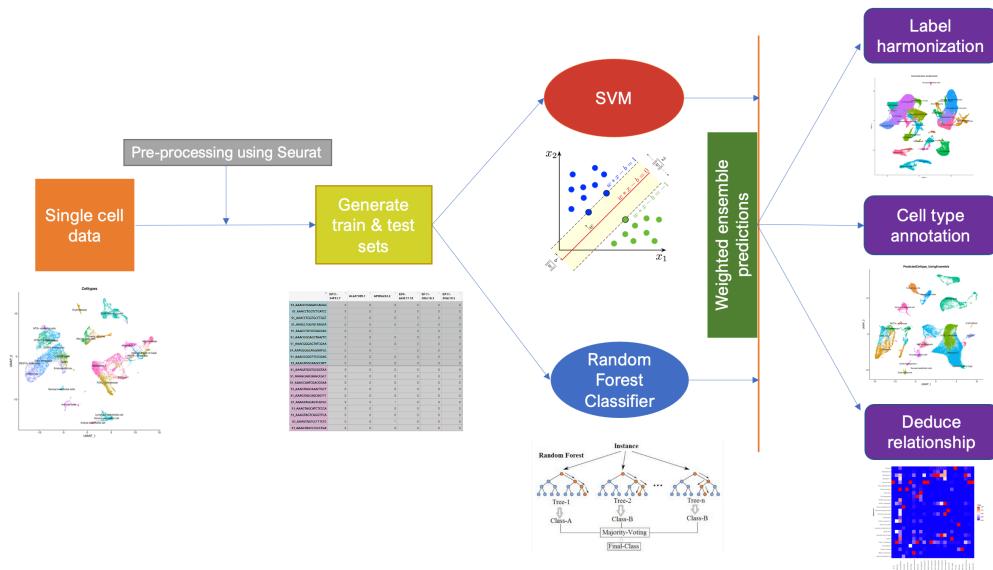
<sup>39</sup> ELeFHAnt has been tested on multiple public datasets involving multiple time points in fetal <sup>40</sup> development, where it was still able to identify cell types. It has also been used internally within

<sup>41</sup> Cincinnati Children's Hospital Medical Center across different projects, including annotation  
<sup>42</sup> of snRNA (single nucleus RNA) sequencing with a harmonized gut cell atlas.

## <sup>43</sup> Overview and Examples

<sup>44</sup> ELeFHAnt makes use of Seurat for the initial input data and pre-processing. It will then  
<sup>45</sup> generate the training and test sets from the reference and query respectively, with optional  
<sup>46</sup> subsampling. SVM and Random Forest are the classifiers that can be used separately or in  
<sup>47</sup> an ensemble. Classification accuracy of both are used to assign weights to the predictions  
<sup>48</sup> from each classifier. These weighted confusion matrices are normalized based on the largest  
<sup>49</sup> number of cells shared among celltypes and assigned clusters. They are then added together  
<sup>50</sup> for the final ensemble predictions.

<sup>51</sup> Figure 1



<sup>52</sup>  
<sup>53</sup> The attributes of our three example datasets of early gut development are shown below, as well  
<sup>54</sup> as those of the integrated dataset. “Fetal” refers to a subset of terminal ileum (TI) data from  
<sup>55</sup> an atlas for human fetal intestinal development called “STAR-FINDer” (Fawknier-Corbett et  
<sup>56</sup> al., 2021). “Gut” refers to a subset of duodenum cell data from the Gut Cell Atlas, which also  
<sup>57</sup> examines intestinal development from 6–10 weeks post-conception (Elmentait et al., 2020).  
<sup>58</sup> Lastly, “Spence” refers a subset of fetal intestinal data from a multi-endodermal organ atlas  
<sup>59</sup> (Yu et al., 2021).

<sup>60</sup> Table 1

	Fetal	Gut	Spence	Integrated
# of Cells	24787	21592	77672	124245
# of Genes	33538	33694	26879	2000
# of Clusters	25	31	26	30
# Cell Types	95	26	24	145
Median nFeature	2029	2858	1913	2121
Median nCount	5141	9988	3754	4833

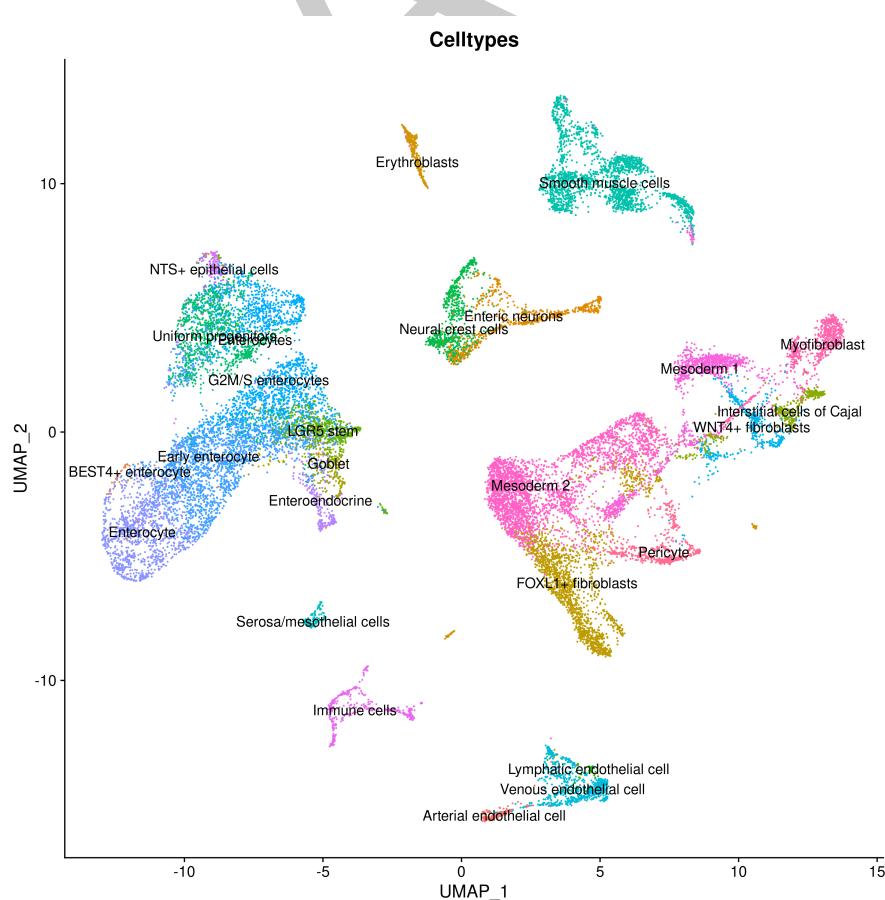
61

## Celltype Annotation

To demonstrate Celltype Annotation using ELeFHAnt we used Gut Cell Atlas ([Elmentait et al., 2020](#)) as reference and Fetal intestinal data ([Yu et al., 2021](#)) as query. Reference and query were downsampled to 200 cells per Celltypes and seurat\_clusters respectively to enable fast computation.

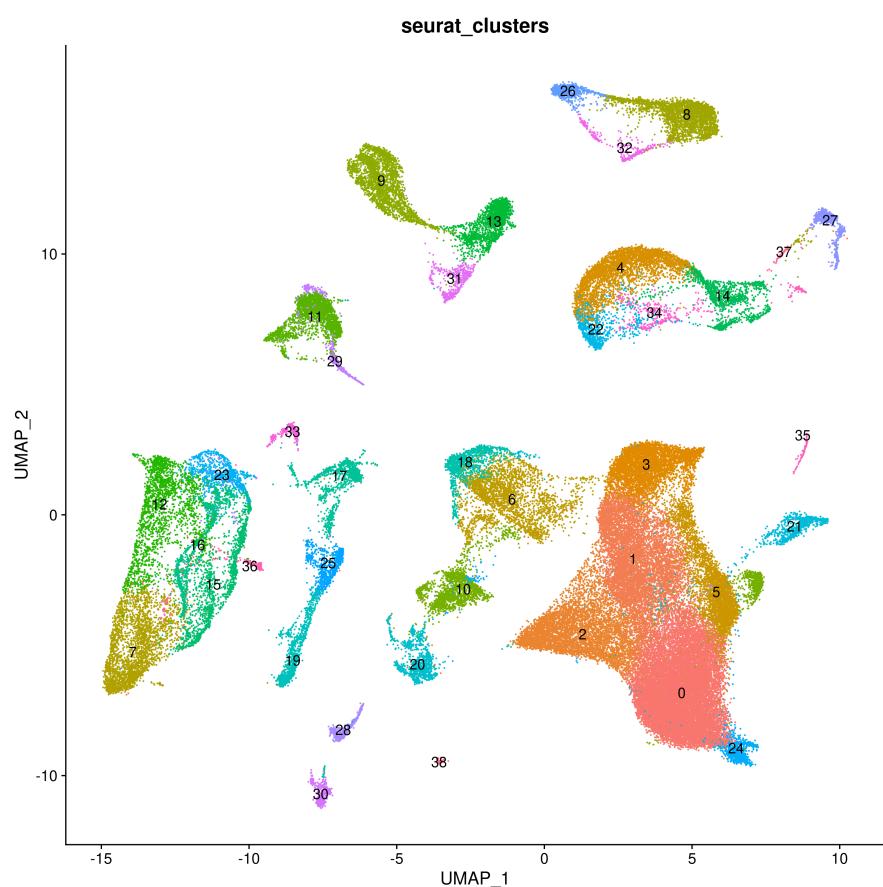
Figure 2

68 A



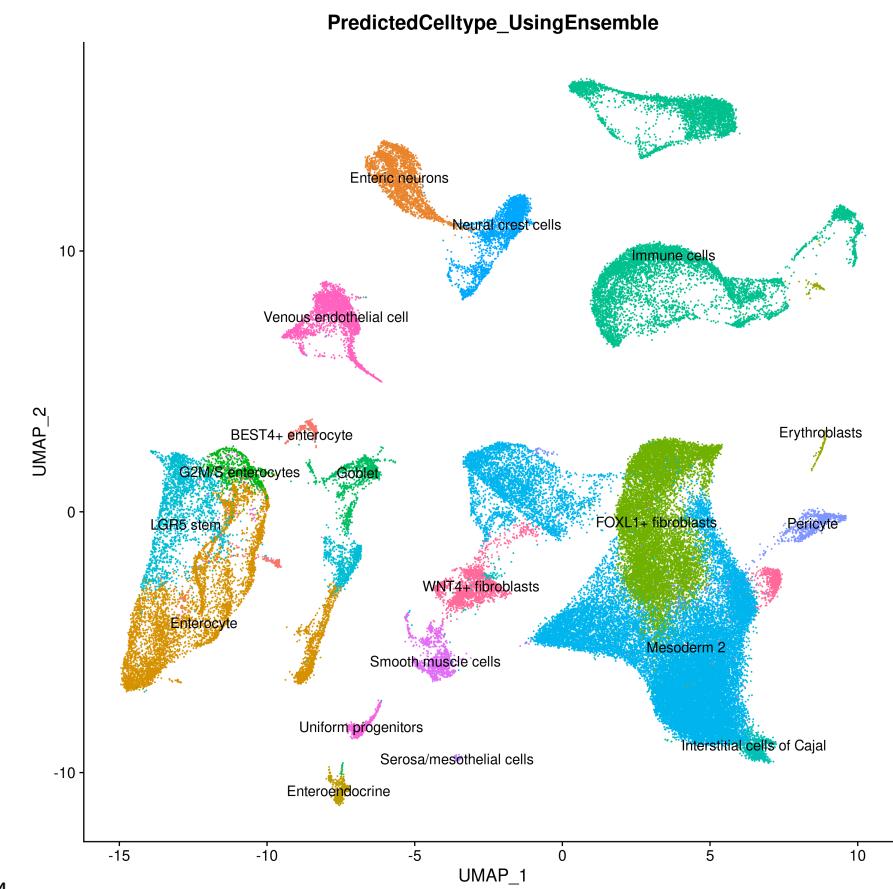
69

70 B



71

72 C

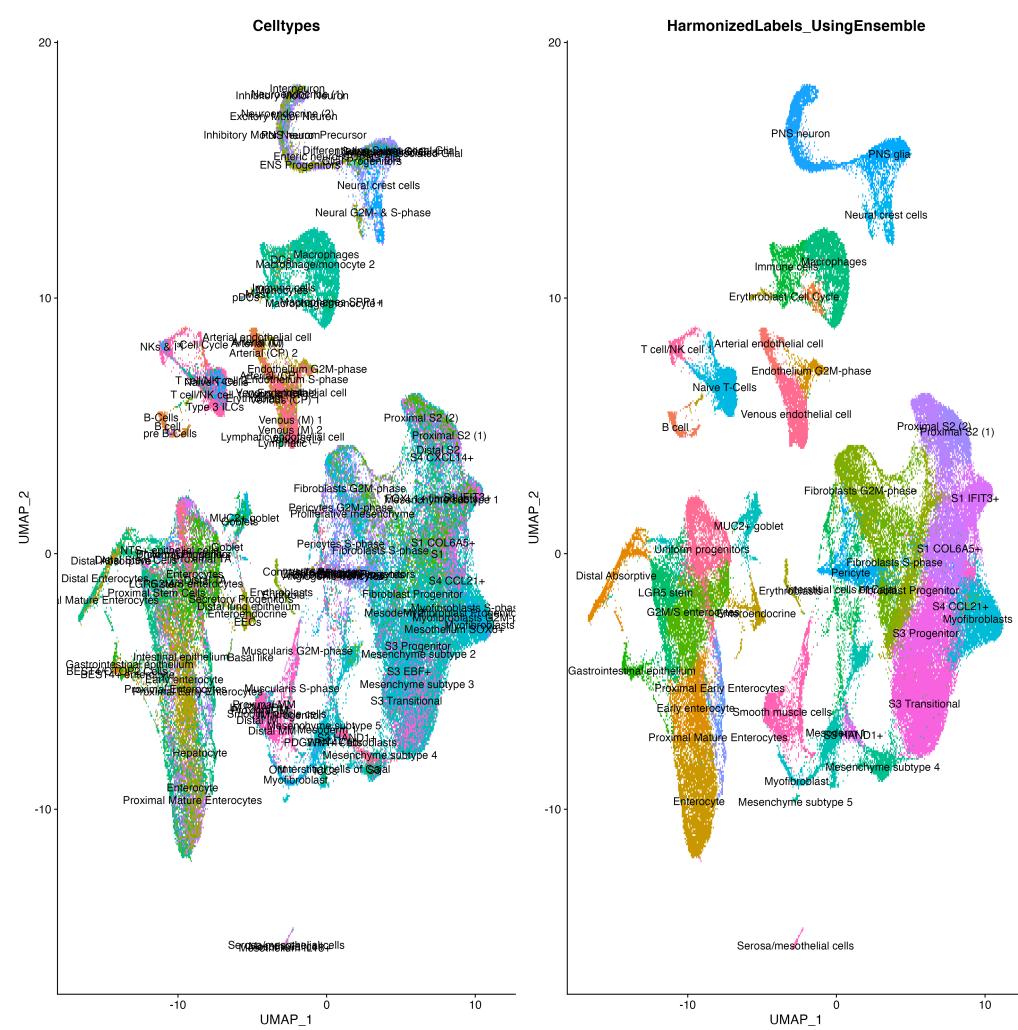


73      4  
 74      (A) represents celltypes in the reference dataset displayed on a UMAP. (B) represents Seurat  
 75      clusters displayed on the query, and (C) represents the predicted celltypes as determined  
 76      by ELeFHAnt's ensemble approach.

#### 77      Label Harmonization

78      To demonstrate LabelHarmonization we used three datasets: 1) Gut Cell Atlas ([Elmentait et al., 2020](#))  
 79      2) Fetal intestinal data ([Yu et al., 2021](#)) from Dr. Spence's Lab 3) Fetal intestine  
 80      data from STAR-FINDer ([Fawknerr-Corbett et al., 2021](#)). Data shown below is based on  
 81      subsetting 200 cells per celltype in each dataset to harmonize the atlas of ~125k cells.

82      Figure 3

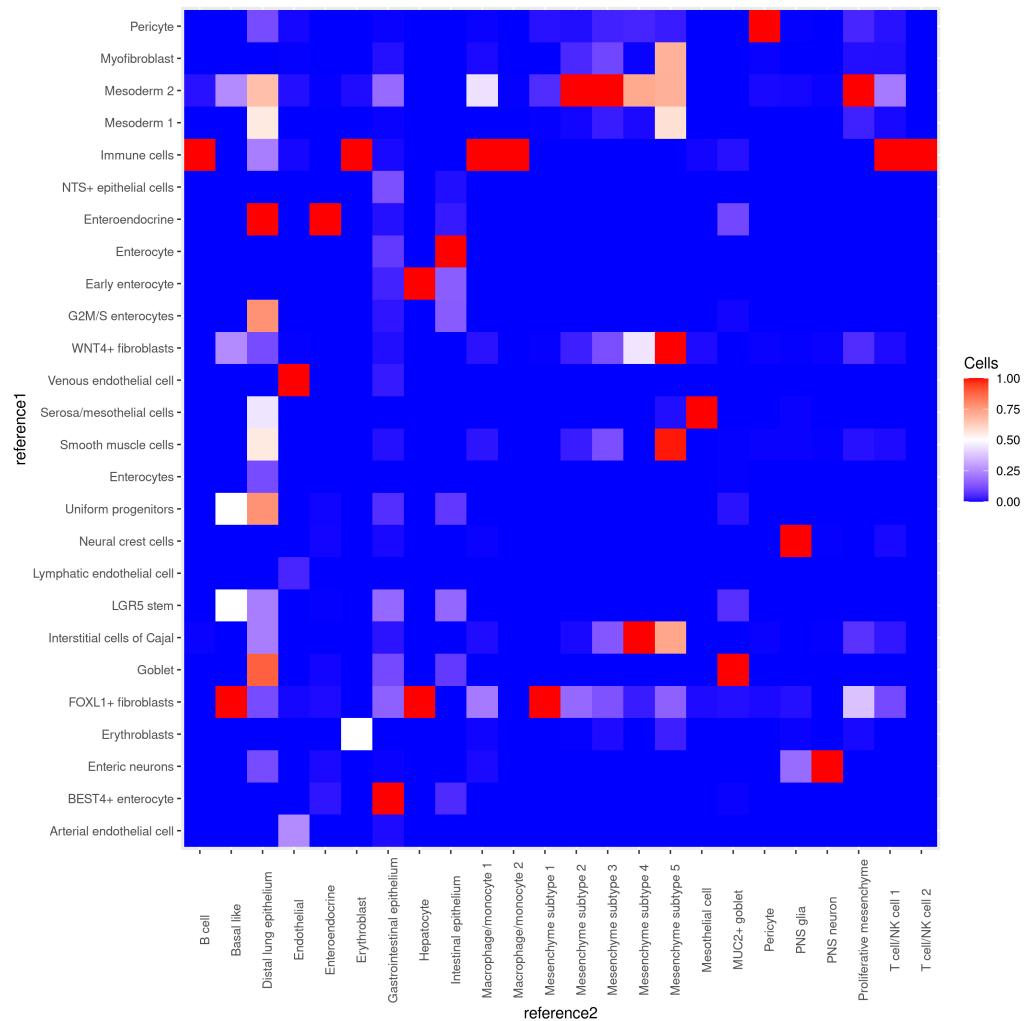


<sup>84</sup> The UMAP on the left represents the 144 total celltypes from the three datasets after integration.  
<sup>85</sup> The UMAP on the right is the result of ELeFHAnt's harmonization with the ensemble  
<sup>86</sup> method, showing the resulting labels for each cluster.

## Deduce Relationship

To demonstrate Deduce Relationship we used two datasets that were also used in the harmonization example: 1) Gut Cell Atlas ([Elmentait et al., 2020](#)) 2) Fetal intestinal data ([Fawkner-Corbett et al., 2021](#)) from Dr. Spence's Lab. Data shown below is based on sub-setting to 500 cells per celltype in each dataset.

92 Figure 4



93  
 94 The heatmap depicts the relationship between celltypes for the two references, with each red  
 95 square showing which cell type in reference 2 best matches a particular celltype in reference  
 96 1.  
 97

## Acknowledgements

98 We would like to thank Drs. Emily Miraldi, Nathan Salomonis, Anil Jegga, and Aaron Zorn  
 99 at Cincinnati Children's Hospital Medical Center for their valuable feedback.

## References

- 101 Abdelaal, T., Michielsen, L., Cats, D., Hoogduin, H., D.and Mei, & Reinders, A., M.and  
 102 Mahfouz. (2019). A comparison of automatic cell identification methods for single-cell  
 103 RNA sequencing data. *Genome Biology*, 20. <https://doi.org/10.1186/s13059-019-1795-z>
- 104 Elmentait, R., Ross, A., Roberts, K., James, K., Ortmann, D., Gomes, T., Nayak, K., Tuck,  
 105 L., Pritchard, S., Bayraktar, O., Heuschkel, R., Vallier, L., Teichmann, S., & Zilbauer, M.  
 106 (2020). Single-Cell Sequencing of Developing Human Gut Reveals Transcriptional Links

- 107 to Childhood Crohn's Disease. *Developmental Cell*, 55. <https://doi.org/10.1016/j.devcel.2020.11.010>
- 108
- 109 Fawkner-Corbett, D., Antanaviciute, A., Parikh, K., Jagielowicz, M., Sousa Gerós, A., Gupta, T., Ashley, N., Khamis, D., Fowler, D., Morrissey, E., Cunningham, C., Johnson, P., Koohy, H., & Simmons, A. (2021). Spatiotemporal analysis of human intestinal development at single-cell resolution. *Cell*, 184. <https://doi.org/10.1016/j.cell.2020.12.016>
- 110
- 111
- 112
- 113 Pasquini, G., Rojo Arias, J. E., Schäfer, P., & Busskamp, V. (2021). Automated methods for cell type annotation on scRNA-seq data. *Computational and Structural Biotechnology Journal*, 19. <https://doi.org/10.1016/j.csbj.2021.01.015>
- 114
- 115
- 116 Stuart, T., Butler, A., Hoffman, P., Hafemeister, C., Papalexi, E., Mauck, W., Hao, Y., Stoeckius, M., Smibert, P., & Satija, R. (2019). Comprehensive Integration of Single-Cell Data. *Cell*, 177. <https://doi.org/10.1016/j.cell.2019.05.031>
- 117
- 118
- 119 Xu, C., Lopez, R., Mehlman, E., Regier, J., Jordan, M. I., & Yosef, N. (2021). Probabilistic harmonization and annotation of single-cell transcriptomics data with deep generative models. *Molecular Systems Biology*, 17. <https://doi.org/10.15252/msb.20209620>
- 120
- 121
- 122 Yu, Q., Kilik, U., Holloway, E., Tsai, Y., Harmel, C., Wu, A., Wu, J., Czerwinski, M., Childs, C., He, Z., Capeling, M., Huang, S., Glass, I., Higgins, P., Treutlein, B., Spence, J., & Camp, J. (2021). Charting human development using a multi-endodermal organ atlas and organoid models. *Cell*, 184. <https://doi.org/10.1016/j.cell.2021.04.028>
- 123
- 124
- 125

DRAFT