

# 1 **SCelVis: Powerful explorative single cell data analysis on** 2 **the desktop and in the cloud**

3  
4 Benedikt Obermayer<sup>1,2,\*</sup>, Manuel Holtgrewe<sup>1,2,\*</sup>, Mikko Nieminen<sup>1,3</sup>, Clemens Messerschmidt<sup>1,2</sup>,  
5 Dieter Beule<sup>1,3</sup>

6  
7 <sup>1</sup> Core Unit Bioinformatics, Berlin Institute of Health, Berlin, Germany

8 <sup>2</sup> Charité – University Medicine Berlin, Berlin, Germany

9 <sup>3</sup> Max Delbrück Center for Molecular Medicine, Berlin, Germany

10 \* these authors contributed equally

11 Corresponding Author:

12 Dieter Beule<sup>1,3</sup>

13 Chariteplatz 1, 10117 Berlin, Germany

14 Email address: [dieter.beule@bihealth.de](mailto:dieter.beule@bihealth.de)

15

## 16 **Abstract**

17 **Background:** Single cell omics technologies present unique opportunities for biomedical and  
18 life sciences from lab to clinic, but the high dimensional nature of such data poses challenges for  
19 computational analysis and interpretation. Furthermore, FAIR data management as well as data  
20 privacy and security become crucial when working with clinical data, especially in cross-  
21 institutional and translational settings. Existing solutions are either bound to the desktop of one  
22 researcher or come with dependencies on vendor-specific technology for cloud storage or user  
23 authentication.

24 **Results:** To facilitate analysis and interpretation of single-cell data by users without  
25 bioinformatics expertise, we present SCelVis, a flexible, interactive and user-friendly app for  
26 web-based visualization of pre-processed single-cell data. Users can survey multiple interactive  
27 visualizations of their single cell expression data and cell annotation, and download raw or  
28 processed data for further offline analysis. SCelVis can be run both on the desktop and cloud  
29 systems, accepts input from local and various remote sources using standard and open protocols,  
30 and allows for hosting data in the cloud and locally.

31 **Methods:** SCelVis is implemented in Python using Dash by Plotly. It is available as a standalone  
32 application as a Python package, via Conda/Bioconda and as a Docker image. All components  
33 are available as open source under the permissive MIT license and are based on open standards  
34 and interfaces, enabling further development and integration with third party pipelines and  
35 analysis components. The GitHub repository is <https://github.com/bihealth/scelvis>.

## 36 Introduction

37 Single-cell omics technologies, in particular single-cell RNA sequencing (scRNA-seq), allow for  
38 the high-throughput profiling of gene expression in thousands to millions of cells with  
39 unprecedented resolution. Recent large-scale efforts to catalogue and describe all human cell  
40 types (Regev et al., 2017) dovetail with ongoing investigations to study cells and tissues in health  
41 and disease, e.g., as proposed by the LifeTime consortium (<https://lifetime-fetflagship.eu>).  
42 Single-cell sequencing could therefore become a routine tool in the clinic for comprehensive  
43 assessments of molecular and physiological alterations in diseased organs as well as systemic  
44 responses, e.g., of the immune system. The enormous scale and high-dimensional nature of the  
45 resulting data presents an ongoing challenge for computational analysis (Stegle, Teichmann, &  
46 Marioni, 2015). Ever more sophisticated methods combining more conventional genomics  
47 approaches with deep learning frameworks (Eraslan, Avsec, Gagneur, & Theis, 2019) allow to  
48 overcome technical limitations and biases and extract multiple layers of information, e.g. from  
49 cell types to lineages and differentiation programs. Many of these methods, their mathematical  
50 background, and the underlying assumptions will remain opaque to users without specific  
51 bioinformatics expertise. At the same time, an in-depth understanding of cell types, their  
52 functional specialization and modification by diseases, and underlying molecular correlates is  
53 often beyond the biological know-how of typical bioinformatics researchers. More than ever,  
54 single-cell omics requires close communication and close collaboration from wet and dry lab  
55 experts. Due to the large amount of data, communication need to be based on interactive  
56 channels (e.g., web-based apps) rather than static tables. Further, as single-cell omics moves  
57 towards the clinic, FAIR (Wilkinson et al., 2016) data management, data privacy, and data  
58 security issues need to be handled appropriately. All employed methods should be able to scale  
59 towards handling a large number of users and even larger numbers of samples.

60 **State of the Art.** Web apps have been used extensively in the single-cell literature and are most  
61 commonly built on Shiny (RStudio Inc., 2014). However, standalone and general-purpose tools  
62 are to our knowledge quite rare. Pagoda (Fan et al., 2016) comes with a simple intuitive web app,  
63 which is limited to Pagoda output and requires manual preprocessing. Cerebro (Hillje, Pelicci, &  
64 Luzi, 2019) is a Shiny web app combined with a Docker container and an Electron  
65 (<https://github.com/electron/electron>) standalone app and provides relatively rich functionality  
66 such as gene set enrichments and quality control statistics, but relies on extensive manual  
67 preprocessing and is not (yet) ready for larger frameworks. On the other hand, the Broad Single  
68 Cell Portal ([https://portals.broadinstitute.org/single\\_cell](https://portals.broadinstitute.org/single_cell)) provides a large-scale web service for a  
69 large number of users and studies. It includes a 10X Genomics data processing pipeline and user  
70 authentication/account management. However, the underlying Docker image strongly depend on  
71 vendor-specific cloud systems such as Google cloud and Broad Firecloud services. Its  
72 implementation thus poses practical hurdles, in particular if it is to be integrated into existing  
73 clinical infrastructure.

## 74 **Materials & Methods**

75 SCelVis is based on Dash by Plotly (Plotly Technologies Inc., 2015) and accepts data in HDF5  
76 format as AnnData objects, which can be created using Scanpy (Wolf, Angerer, & Theis, 2018).  
77 It also provides conversion functionality from raw text or 10X Genomics CellRanger output. The  
78 built-in converter is accessible from the command line and a web-based user interface (Figure 1).  
79 It allows for converting pipeline output with an optional description file into a single AnnData  
80 HDF5 file. One HDF5 file or a folder containing multiple such files can then be provided to  
81 SCelVis for visualization, and data sets can be selected for exploration on the graphic web  
82 interface. To enable both local and cloud access, data can be read from the file system or remote  
83 data sources via the standard internet protocols FTP, SFTP, and HTTP(S). SCelVis also provides  
84 data access through the open source iRODS protocol (Rajasekar et al., 2010) or the widely-used  
85 Amazon S3 object storage protocol. The data sources can be given on the command line and as  
86 environment variables as is best practice for cloud deployments (Adam Wiggins, 2011). The  
87 latter allows for easy “serverless” and cloud deployments.

88 SCelVis is built around two perspectives on single-cell data (Figure 1). On the one hand, it  
89 provides a cell-based view, where users can browse and investigate cell annotations (such as cell  
90 type) and cell-specific statistics (such as sequencing depth or cell type proportions) in multiple  
91 visualizations, e.g., on a t-SNE or UMAP embedding, as violin plots or bar charts. On the other  
92 hand, it provides a gene-based view that lets users explore gene expression in multiple  
93 visualizations on embeddings or as violin or dot plots. Relevant genes can be specified by hand  
94 or selected directly from lists of marker genes.

95 The source code is available under the permissive MIT license on the GitHub repository at  
96 <https://github.com/bihealth/scelvis>. The software can be run both in the cloud and on workstation  
97 desktops via Docker.

## 98 **Usage Example**

99 We provide two example datasets within our Github repository (see above, it also contains a link  
100 to a public demonstration instance). First, a small synthetic simulated dataset created for  
101 illustration purposes, and secondly a publicly available processed scRNA-seq dataset from 10X  
102 Genomics containing ~1000 cells of a mix of human HEK293T and murine NIH3T3 cells  
103 (Figure 2).

## 104 **Discussion**

105 In this manuscript, we have presented SCelVis, a method for the interactive visualization of  
106 single-cell RNA-seq data. It provides easy-to-use yet flexible means of scRNA-seq data  
107 exploration for researchers without computational background. SCelVis takes processed data,  
108 e.g., provided by CellRanger or a bioinformatics collaboration partner, as input, and focuses  
109 solely on visualization and explorative analysis. Great care has been taken to make the method  
110 flexible in usage and deployment. It can be used both on a researcher’s desktop with minimal  
111 training yet its usage scales up to a cloud deployment. Data can be read from local file systems  
112 but also from a variety of remote data sources, e.g., via the widely deployed (S)FTP, S3, and

113 HTTP(S) protocols. This allows for deploying it in a Docker container on “serverless” cloud  
114 systems. As both the application and data can be hosted on the network or cloud systems, the  
115 application facilitates cross-institutional research. For example, a sequencing or bioinformatics  
116 core unit can use it for giving access to non-computational collaboration partners over the  
117 internet. This is particularly interesting as it comes with no dependency on any vendor-specific  
118 technology such as the Google or Facebook authentication that appears to become pervasive in  
119 today’s life science.

## 120 **Conclusions**

121 SCelVis is a flexible and powerful method for the visualization of single-cell RNA-seq  
122 experiments and the explorative data analysis thereof. It comes with a number of unique features,  
123 in particular complete independence of vendor-specific software or services. At the same time, it  
124 remains simple enough to be integrated as a component in more complex framework.

## 125 **Acknowledgements**

126 The example dataset for the 1:1 mixture of human and mouse cells processed with CellRanger  
127 (v3) was taken from the 10X genomics website [https://support.10xgenomics.com/single-cell-gene-expression/datasets/3.0.0/hgmm\\_1k\\_v3](https://support.10xgenomics.com/single-cell-gene-expression/datasets/3.0.0/hgmm_1k_v3).

## 129 **References**

- 130 Adam Wiggins. (2011). The Twelve-Factor App. Retrieved March 22, 2019, from  
131 <https://12factor.net/>
- 132 Eraslan, G., Avsec, Ž., Gagneur, J., & Theis, F. J. (2019). Deep learning: new computational  
133 modelling techniques for genomics. *Nature Reviews Genetics*, 20(7), 389–403.  
134 <https://doi.org/10.1038/s41576-019-0122-6>
- 135 Fan, J., Salathia, N., Liu, R., Kaeser, G. E., Yung, Y. C., Herman, J. L., ... Kharchenko, P. V.  
136 (2016). Characterizing transcriptional heterogeneity through pathway and gene set  
137 overdispersion analysis. *Nature Methods*, 13(3), 241–244.  
138 <https://doi.org/10.1038/nmeth.3734>
- 139 Hillje, R., Pelicci, P. G., & Luzi, L. (2019). Cerebro: Interactive visualization of scRNA-seq  
140 data. *BioRxiv*, 631705. <https://doi.org/10.1101/631705>
- 141 Plotly Technologies Inc. (2015). *Collaborative data science*. Retrieved from <https://plot.ly>
- 142 Rajasekar, A., Moore, R., Hou, C.-Y., Lee, C. A., Marciano, R., de Torcy, A., ... Zhu, B. (2010).  
143 iRODS Primer: Integrated Rule-Oriented Data System. *Synthesis Lectures on Information*  
144 *Concepts, Retrieval, and Services*, 2(1), 1–143.  
145 <https://doi.org/10.2200/S00233ED1V01Y200912ICR012>
- 146 Regev, A., Teichmann, S. A., Lander, E. S., Amit, I., Benoist, C., Birney, E., ... Human Cell  
147 Atlas Meeting Participants. (2017). The Human Cell Atlas. *ELife*, 6.  
148 <https://doi.org/10.7554/eLife.27041>
- 149 RStudio Inc. (2014). *shiny: Easy web applications in R*. Retrieved from <http://shiny.rstudio.com>
- 150 Stegle, O., Teichmann, S. A., & Marioni, J. C. (2015). Computational and analytical challenges  
151 in single-cell transcriptomics. *Nature Reviews Genetics*, 16(3), 133–145.  
152 <https://doi.org/10.1038/nrg3833>
- 153 Wilkinson, M. D., Dumontier, M., Aalbersberg, Ij. J., Appleton, G., Axton, M., Baak, A., ...

- 154 Mons, B. (2016). The FAIR Guiding Principles for scientific data management and  
155 stewardship. *Scientific Data* 2016 3.
- 156 Wolf, F. A., Angerer, P., & Theis, F. J. (2018). SCANPY: large-scale single-cell gene expression  
157 data analysis. *Genome Biology*, 19(1), 15. <https://doi.org/10.1186/s13059-017-1382-0>  
158

hdf5 files, local or remote source:

- (s)ftp / http(s)
- iRODs (ssl / PAM authentication / tickets)
- S3

other formats

- raw text
- cellranger output

SCelVis  
local or cloud instance

scelvis convert

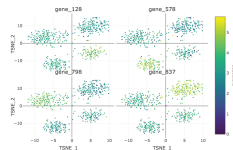
cell-centric view

gene-centric view

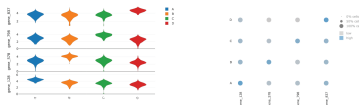
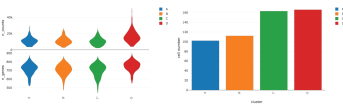
clustering /  
cell annotation



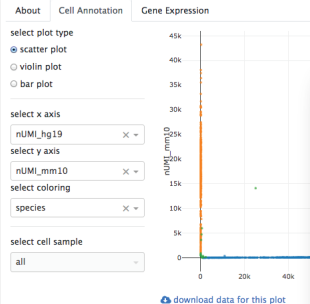
expression /  
marker genes



QC /  
proportions



**A:** scatterplots for cell annotation (e.g., clustering) or quality control statistics



select plot type

 scatter plot violin plot dot plot

select x axis

TSNE1

select y axis

TSNE2

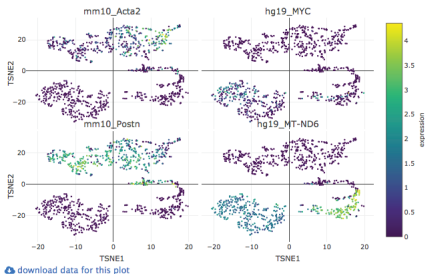
select gene(s)

 mm10\_Acta2 hg19\_MYC mm10\_Postn hg19\_MT-ND6 use marker table for selection

select cell sample

all

use selected genes



	Cluster	gene	log2_fc	mean_counts	p
<input checked="" type="checkbox"/>	Cluster_1	mm10_Acta2	5.397	17.071	0
<input type="checkbox"/>	Cluster_1	mm10_Tagln	4.227	4.673	0
<input type="checkbox"/>	Cluster_1	mm10_Txn1	3.269	15.562	0
<input type="checkbox"/>	Cluster_1	mm10_Ctst1	3.268	11.554	0
<input type="checkbox"/>	Cluster_1	mm10_Rps2	3.143	59.141	0
<input type="checkbox"/>	Cluster_1	mm10_Drb1	3.564	4.067	0

**B:** plot gene expression as scatter, violin or dot plots (select genes from table)