





scmorph: a Python package for analysing single-cell morphological profiles

Jesko Wagner¹, Hugh Warden¹, Ava Khamseh^{1,2,3}, and Sjoerd Viktor Beentjes^{1,3}

¹ MRC Human Genetics Unit, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh EH4 2XU, UK ² School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, UK ³ School of Mathematics and Maxwell Institute for Mathematical Sciences, University of Edinburgh, Edinburgh EH9 3FD, UK  Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

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

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

scmorph is a Python package to analyse single-cell data from morphological profiling experiments which generate large tabular data. scmorph combines domain-specific methods such as single-cell hit calling and batch correction with the versatile and scalable [scverse](#) tools to offer feature selection, dimensionality reduction and more. Overall, scmorph brings together a host of single-cell morphological profiling methods, making it applicable for a wide range of experimental designs and workflows.

Statement of need

Morphological profiling has become an essential tool in biology and drug discovery, but there is a lack of open-source software for analysing single-cell morphological data. Existing solutions are commercial, do not scale to large datasets, or do not offer single-cell specific methods ([Omta et al., 2016](#); [Serrano et al., 2025](#)). scmorph offers a comprehensive set of methods for analysing single-cell morphological data, which do not require averaging of features across cells. By skipping the averaging or profiles, scmorph retains the heterogeneity of cell populations. This enables hit-calling that is sensitive to changes in subpopulations and recapitulating dynamic processes such as differentiation. By integrating with the growing scverse of single-cell tools, scmorph also opens up advanced processing capabilities including access to deep learning tools ([Wolf et al., 2018](#)).

Briefly, scmorph provides five modules to analyze morphological profiles:

- Reading and writing (IO). scmorph allows reading data from csv, sql, sqlite, and h5ad files, including from the popular CellProfiler software ([Stirling et al., 2021](#)). Once converted, scmorph works with AnnData objects stored as h5ad, which track processing steps and can easily be written to disk ([Virshup et al., 2024](#)).
- Quality control. scmorph integrates two levels of unsupervised quality control: image-level and single-cell level. Image-level correction is performed with a kNN-based outlier detection method, whereas single-cell profiles that are outliers are detected via pyod ([Li et al., 2022](#)).
- Preprocessing. Provided functions perform feature selection, compute PCA coordinates, and optionally aggregate data. For the first time in the field of morphological profiling, scmorph integrates scone as batch correction function, which retains interpretability of features ([Cole et al., 2019](#)). Additionally, the integrated feature selection methods can remove features associated with known confounders or with high correlation structures, as is common in morphological profiling experiments ([Kruskal & Wallis, 1952](#); [Lin &](#)

- 42 [Han, 2021](#)).
- 43 ■ Plotting. scmorph uses scanpy for easy plotting of PCA and UMAP coordinates, either
- 44 in 2D or as cumulative densities, which can be useful for identifying technical artifacts
- 45 such as batch effects ([Wolf et al., 2018](#)). It also provides methods for plotting features
- 46 per experimental group, such as plates.
- 47 ■ Downstream analysis. For experiments focused on profiling non-dynamic responses, such
- 48 as a small molecule library, scmorph integrates functions to perform hit calling from
- 49 single-cell profiles. Specifically, during hit calling scmorph embeds single-cell profiles into
- 50 PCA space and computes the Mahalanobis distance of cells to the medoid of untreated
- 51 control cells. It then compares the distances of treated cells to those of untreated cells
- 52 by use of the Kolmogorov–Smirnov statistic. For dynamic systems such as differentiating
- 53 cells, scmorph incorporates differential trajectory inference modelling via slingshot and
- 54 condiments through the rpy2 translation layer ([Roux de Bézieux et al., 2024](#); [Street et](#)
- 55 [al., 2018](#)).

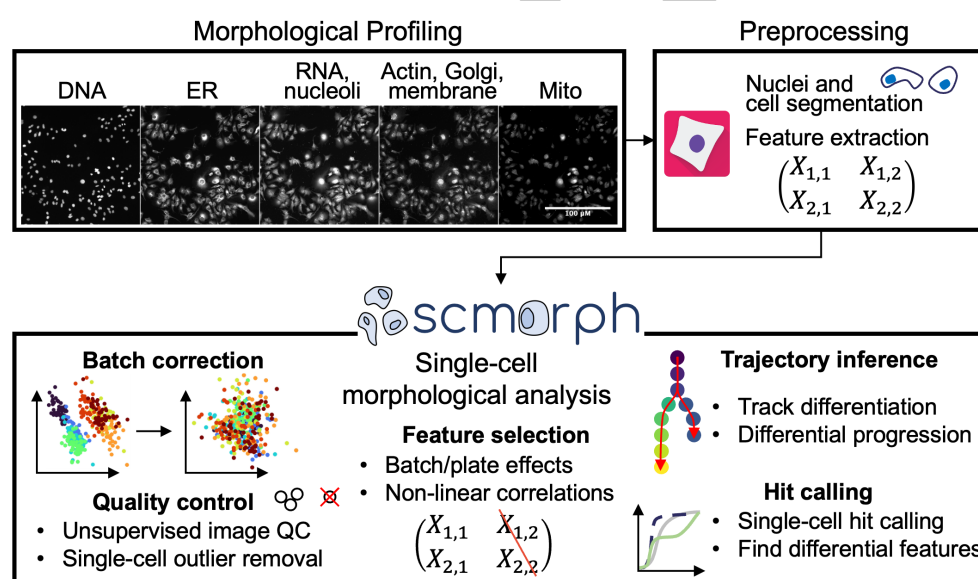


Figure 1: Overview of scmorph functionality. scmorph processes profiles generated with software such as CellProfiler to facilitate downstream analysis by performing batch correction, image- and single-cell QC, feature selection, hit calling and trajectory inference. All methods are built with single-cell analysis in mind and do not require subsampling.

56 In contrast to the commonly used pycytominer package ([Serrano et al., 2025](#)) and SPACe

57 ([Stossi et al., 2024](#)), scmorph offers (i) interpretable batch correction techniques compatible

58 with single-cell profiles, (ii) enhanced feature selection with an adapted Chatterjee correlation

59 coefficient or Kruskal-Wallis test ([Kruskal & Wallis, 1952](#); [Lin & Han, 2021](#)), and (iii)

60 lineage trajectory inference ([Roux de Bézieux et al., 2024](#); [Street et al., 2018](#)). Compared to

61 pycytominer, scmorph also performs single-cell based hit calling. And unlike SPACe, scmorph

62 is agnostic to the segmentation and feature extraction methods used upstream and therefore

63 compatible with CellProfiler. scmorph also benefits from improvements of AnnData and scanpy,

64 such as enabling out-of-core processing crucial to big data analysis ([Virshup et al., 2024](#); [Wolf](#)

65 [et al., 2018](#)).

66 scmorph aims to complement rather than replace traditional aggregate analysis. We recommend

67 its use in experiments where heterogeneous cells or cell-state specific responses to treatments are

68 expected. In more traditional drug discovery campaigns, scmorph may be of more limited use

69 due to slower runtimes and because outputs require more careful interpretation. Nevertheless,

we envision that tools offered in scmorph such as the interpretable batch correction could find broad application prior to aggregation even in traditional imaging campaigns.

Already, scmorph has been used to quality control morphological profiling experiments involving differentiating liver cells (Graham et al., 2025). scmorph is also involved in three projects involving small compound and microRNA perturbations in the domains of drug discovery and fundamental research, spanning datasets of >20M cells. Going forward, we envision that scmorph will enable analysis of complex and large morphological profiling experiments.

Acknowledgements

JW and HW are supported by an MRC PhD studentship (grant no. MC_ST_00035).

References

- Cole, M. B., Risso, D., Wagner, A., DeTomaso, D., Ngai, J., Purdom, E., Dudoit, S., & Yosef, N. (2019). Performance Assessment and Selection of Normalization Procedures for Single-Cell RNA-Seq. *Cell Systems*, 8(4), 315–328.e8. <https://doi.org/10.1016/j.cels.2019.03.010>
- Graham, R. E., Zheng, R., Wagner, J., Unciti-Broceta, A., Hay, D. C., Forbes, S. J., Gadd, V. L., & Carragher, N. O. (2025). Single-cell morphological tracking of cell states to identify small-molecule modulators of liver differentiation. *iScience*, 28(2), 111871. <https://doi.org/10.1016/j.isci.2025.111871>
- Kruskal, W. H., & Wallis, W. A. (1952). Use of Ranks in One-Criterion Variance Analysis. *Journal of the American Statistical Association*, 47(260), 583–621. <https://doi.org/10.1080/01621459.1952.10483441>
- Li, Z., Zhao, Y., Hu, X., Botta, N., Ionescu, C., & Chen, G. H. (2022). ECOD: Unsupervised Outlier Detection Using Empirical Cumulative Distribution Functions. *IEEE Transactions on Knowledge and Data Engineering*, 1–1. <https://doi.org/10.1109/tkde.2022.3159580>
- Lin, Z., & Han, F. (2021). On boosting the power of Chatterjee's rank correlation. *arXiv*. <https://doi.org/10.48550/arxiv.2108.06828>
- Omta, W. A., van Heesbeen, R. G., Pagliero, R. J., van der Velden, L. M., & Lelieveld, D. (2016). HC StratoMineR: A web-based tool for the rapid analysis of high-content datasets. *ASSAY and Drug Development Technologies*, 14(8), 439–452. <https://doi.org/gt2g9c>
- Roux de Bézieux, H., Van den Berge, K., Street, K., & Dudoit, S. (2024). Trajectory inference across multiple conditions with condiments. *Nature Communications*, 15(1), 833. <https://doi.org/10.1038/s41467-024-44823-0>
- Serrano, E., Chandrasekaran, S. N., Bunten, D., Brewer, K. I., Tomkinson, J., Kern, R., Bornholdt, M., Fleming, S. J., Pei, R., Arevalo, J., Tsang, H., Rubinetti, V., Tromans-Coia, C., Becker, T., Weisbart, E., Bunne, C., Kalinin, A. A., Senft, R., Taylor, S. J., ... Way, G. P. (2025). Reproducible image-based profiling with Pycytominer. *Nature Methods*. <https://doi.org/10.1038/s41592-025-02611-8>
- Stirling, D. R., Swain-Bowden, M. J., Lucas, A. M., Carpenter, A. E., Cimini, B. A., & Goodman, A. (2021). CellProfiler 4: Improvements in speed, utility and usability. *BMC Bioinformatics*, 22(1), 433. <https://doi.org/10.1186/s12859-021-04344-9>
- Stossi, F., Singh, P. K., Marini, M., Safari, K., Szafran, A. T., Rivera Tostado, A., Candler, C. D., Mancini, M. G., Mosa, E. A., Bolt, M. J., Labate, D., & Mancini, M. A. (2024). SPACE: An open-source, single-cell analysis of Cell Painting data. *Nature Communications*, 15(1), 10170. <https://doi.org/10.1038/s41467-024-54264-4>
- Street, K., Risso, D., Fletcher, R. B., Das, D., Ngai, J., Yosef, N., Purdom, E., & Dudoit, S.

- 114 (2018). Slingshot: Cell lineage and pseudotime inference for single-cell transcriptomics.
115 *BMC Genomics*, 19(1), 477. <https://doi.org/10.1186/s12864-018-4772-0>
- 116 Virshup, I., Rybakov, S., Theis, F. J., Angerer, P., & Wolf, F. A. (2024). Anndata: Access
117 and store annotated data matrices. *Journal of Open Source Software*, 9(101), 4371.
118 <https://doi.org/10.21105/joss.04371>
- 119 Wolf, F. A., Angerer, P., & Theis, F. J. (2018). SCANPY: Large-scale single-cell gene expression
120 data analysis. *Genome Biology*, 19(1), 15. <https://doi.org/10.1186/s13059-017-1382-0>

DRAFT