

SpatialProteomicsNet : A unified interface for spatial omics data access for computer vision and machine learning

2025-05-28

Summary

SpatialProteomicsNet is an open-source Python package that provides a harmonized and standardized interface for accessing spatial proteomics and multiplexed imaging datasets, including imaging mass cytometry (IMC) (Giesen et al. 2014) and multiplexed ion beam imaging (MIBI) (Keren et al. 2019) data. The package enables researchers to load raw spatially-resolved proteomics data from multiple studies in a unified format, apply and retrieve data structures ready for downstream machine learning analysis or model training. By focusing on open-source raw data processing and enforcing common data schemas (e.g., standardized image and single-cell data formats), **SpatialProteomicsNet** promotes reproducible and efficient research in computational and spatial biology. The library is designed to serve the broader community working on spatial proteomics by easing data access and integration into machine learning workflows.

Statement of Need

Spatially-resolved proteomics, recently named Nature Method of the Year 2024 (“Method of the Year 2024: Spatial Proteomics” 2024), enable the quantification of proteins in single cells within their tissue context, revealing intricate aspects of spatial cellular arrangement and communication. In the context of cancer, these advancements provide unprecedented insights into the heterogeneity of the tumor and its microenvironment, and the underlying mechanisms affecting tumor initiation, progression, and response to treatment (Lewis et al. 2021). IMC and MIBI-TOF are among the most popular technologies, with dozens of high-dimensional datasets made publicly available per year. The increasing availability of these datasets has fueled algorithmic development in machine learning and computer vision. Numerous models that perform a variety of tasks, such as cell segmentation (Greenwald et al. 2022), cell type annotation (Geuenich et al. 2021), representation learning (Wenckstern et al. 2025) or heterogeneity analysis (Martinelli and Rapsomaniki 2022) tailored to spatial proteomics data have been recently developed, with corresponding widely used packages.

However, a critical gap hindering model development, reproducibility and cross-study analyses is the lack of unified frameworks to access and process the data. Spatial proteomics datasets,

often deposited in public repositories such as Zenodo (European Organization For Nuclear Research and OpenAIRE 2013) or Figshare (“Figshare - Credit for All Your Research” n.d.), typically contain a collection of components, such as raw and preprocessed images, segmentation masks, extracted single-cell intensities, panel descriptions and associated clinical metadata, uploaded in disparate, non-standardized formats (e.g., mixed .tiff, .csv, custom JSONs), with varying metadata structures and inconsistent preprocessing that vary greatly between studies and labs. Working with these fragmented datasets implies a significant time investment for researchers to locate and download the data, and write custom scripts to handle their specific data structure, creating barriers to entry, complicating usage and hindering robust benchmarking. While existing data frameworks developed by the spatial transcriptomics community such as SpatialData (Marconato et al. 2025) and Pysodb (Yuan et al. 2023) are gaining popularity and can be extended to spatial proteomics, they often come with heavier dependencies and general-purpose abstractions that may be unnecessarily complex for researchers focused on fast, standardized access to real-world IMC or MIBI datasets.

SpatialProteomicsNet is an open-source Python package that addresses these gaps by: Providing a lightweight, unified interface to widely-used curated spatial proteomics datasets. Abstracting dataset-specific structure, letting users access data components (images, masks, metadata) through a consistent schema. Supporting reproducible preprocessing via modular, reusable interfaces for common pipeline steps. Facilitating integration in machine learning and computer vision models by streamlining dataset loading into standard formats. Encouraging community contributions for expanding and maintaining harmonized dataset access. This unified approach allows scientists to abstract away dataset-specific idiosyncrasies and focus on biological and analytical questions rather than data wrangling. **SpatialProteomicsNet** is intentionally minimal, tailored to machine learning and computer vision workflows (e.g., loading images, masks, and cell-level CSVs into memory with minimal setup) without depending on larger ecosystem packages (e.g., anndata, xarray, zarr, dask). **SpatialProteomicsNet** gives immediate access to curated datasets with ready-to-use utilities, eliminating the need to write custom loaders or parse inconsistent formats. As such, it is particularly friendly to the growing community of ML developers, researchers, and engineers entering the emerging field of spatial biology. By harmonizing data access, our package enables more straightforward integration of spatial proteomics data into machine learning and modeling frameworks, ultimately accelerating biomedical discovery.

Supported Datasets

The package supports the following public spatial proteomics datasets:

- **Jackson2020** – IMC of breast cancer (Jackson et al. 2020)
- **Danenberg2022** – IMC of breast cancer (Danenberg et al. 2022)
- **Cords2024** – IMC of Non-small cell lung cancer (Cords et al. 2024)
- **Keren2018** – IMC of triple-negative breast cancer (Keren et al. 2018)

	name	# images	# masks	# markers	# cells	# samples
0	Danenberg2022	794	794	39	1123466	794
1	Cords2024	2070	2070	43	5984454	2072
2	Jackson2020	735	735	35	1224411	735

	name	# images	# masks	# markers	# cells	# samples
3	Keren2018	41	41	36	201656	41

Table 1: Summary statistics of supported spatial proteomics datasets in the package.

Additionally, dummy datasets are provided to mimic real data structure for development and testing purposes.

Each dataset is accessible through a standardized class interface that mimics the pytorch lightning (Falcon and team 2019) philosophy and includes methods for downloading, preparing, and accessing processed components (images, masks, metadata, and spatial coordinates). These datasets follow consistent naming conventions and data schemas, making them immediately usable for downstream tasks.

Example Usage

```
from ai4bmr_datasets import Jackson2020
from pathlib import Path

dataset = Jackson2020(base_dir=Path("<PATH>"))
dataset.prepare_data()
dataset.setup(image_version="published", mask_version="published")

print(dataset.sample_ids)  # list of sample IDs
print(dataset.images)      # list of images
print(dataset.masks)       # list of masks

dataset.setup(image_version="published", mask_version="published",
              feature_version='published', load_intensity=True,
              metadata_version='published', load_metadata=True,
              )
print(dataset.intensity.shape)  # cell x marker matrix
print(dataset.metadata.shape)  # cell x annotation matrix
```

Conclusion

SpatialProteomicsNet lowers the technical barrier to working with spatial proteomics data by providing unified, open access to several published datasets and processing routines. Its modular design and standardized outputs make it a practical tool for researchers developing computational methods in spatial biology. We welcome contributions and extensions from the community and envision this package as a foundation for reproducible spatial proteomics analysis.

References

- Cords, Lena, Stefanie Engler, Martina Haberecker, Jan Hendrik Rüschhoff, Holger Moch, Natalie De Souza, and Bernd Bodenmiller. 2024. "Cancer-Associated Fibroblast Phenotypes Are Associated with Patient Outcome in Non-Small Cell Lung Cancer." *Cancer Cell* 42 (3): 396–412.e5. <https://doi.org/10.1016/j.ccell.2023.12.021>.
- Danenberg, Esther, Helen Bardwell, Vito R. T. Zanutelli, Elena Provenzano, Suet Feung Chin, Oscar M. Rueda, Andrew Green, et al. 2022. "Breast Tumor Microenvironment Structures Are Associated with Genomic Features and Clinical Outcome." *Nature Genetics* 54 (5): 660–69. <https://doi.org/10.1038/s41588-022-01041-y>.
- European Organization For Nuclear Research, and OpenAIRE. 2013. "Zenodo." CERN. <https://doi.org/10.25495/7GXK-RD71>.
- Falcon, William, and The PyTorch Lightning team. 2019. "PyTorch Lightning." <https://doi.org/10.5281/zenodo.3828935>.
- "Figshare - Credit for All Your Research." n.d. <https://figshare.com/>. Accessed June 6, 2025.
- Geuenich, Michael J., Jinyu Hou, Sunyun Lee, Shanza Ayub, Hartland W. Jackson, and Kieran R. Campbell. 2021. "Automated Assignment of Cell Identity from Single-Cell Multiplexed Imaging and Proteomic Data." *Cell Systems* 12 (12): 1173–1186.e5. <https://doi.org/10.1016/j.cels.2021.08.012>.
- Giesen, Charlotte, Hao A. O. Wang, Denis Schapiro, Nevena Zivanovic, Andrea Jacobs, Bodo Hattendorf, Peter J. Schöffler, et al. 2014. "Highly Multiplexed Imaging of Tumor Tissues with Subcellular Resolution by Mass Cytometry." *Nature Methods* 11 (4): 417–22. <https://doi.org/10.1038/nmeth.2869>.
- Greenwald, Noah F., Geneva Miller, Erick Moen, Alex Kong, Adam Kagel, Thomas Dougherty, Christine Camacho Fullaway, et al. 2022. "Whole-Cell Segmentation of Tissue Images with Human-Level Performance Using Large-Scale Data Annotation and Deep Learning." *Nature Biotechnology* 40 (4): 555–65. <https://doi.org/10.1038/s41587-021-01094-0>.
- Jackson, Hartland W., Jana R. Fischer, Vito R. T. Zanutelli, H. Raza Ali, Robert Mechera, Savas D. Soysal, Holger Moch, et al. 2020. "The Single-Cell Pathology Landscape of Breast Cancer." *Nature* 578 (7796): 615–20. <https://doi.org/10.1038/s41586-019-1876-x>.
- Keren, Leeat, Marc Bosse, Steve Thompson, Tyler Risom, Kausalia Vijayaragavan, Erin McCaffrey, Diana Marquez, et al. 2019. "MIBI-TOF: A Multiplexed Imaging Platform Relates Cellular Phenotypes and Tissue Structure." *Science Advances* 5 (10): eaax5851. <https://doi.org/10.1126/sciadv.aax5851>.
- Keren, Leeat, Diana Marquez, Marc Bosse, Roshan Angoshtari, Samir Jain, Sushama Varma, Soo Ryum Yang, et al. 2018. "A Structured Tumor-Immune Microenvironment in Triple Negative Breast Cancer Revealed by Multiplexed Ion Beam Imaging." *Cell* 174 (6): 1373–1387.e19. <https://doi.org/10.1016/j.cell.2018.08.039>.
- Lewis, Sabrina M., Marie-Liesse Asselin-Labat, Quan Nguyen, Jean Berthelet, Xiao Tan, Verena C. Wimmer, Delphine Merino, Kelly L. Rogers, and Shalin H. Naik. 2021. "Spatial Omics and Multiplexed Imaging to Explore Cancer Biology." *Nature Methods* 18 (9): 997–1012. <https://doi.org/10.1038/s41592-021-01203-6>.
- Marconato, Luca, Giovanni Palla, Kevin A. Yamauchi, Isaac Virshup, Elyas Heidari, Tim Treis, Wouter-Michiel Vierdag, et al. 2025. "SpatialData: An Open and Universal Data Framework for Spatial Omics." *Nature Methods* 22 (1): 58–62. <https://doi.org/10.1038/s41592-024-02212-x>.
- Martinelli, Adriano Luca, and Maria Anna Rapsomaniki. 2022. "ATHENA: Analysis of

- Tumor Heterogeneity from Spatial Omics Measurements.” Edited by Hanchuan Peng. *Bioinformatics* 38 (11): 3151–53. <https://doi.org/10.1093/bioinformatics/btac303>.
- “Method of the Year 2024: Spatial Proteomics.” 2024. *Nature Methods* 21 (12): 2195–96. <https://doi.org/10.1038/s41592-024-02565-3>.
- Wenckstern, Johann, Eeshaan Jain, Kiril Vasilev, Matteo Pariset, Andreas Wicki, Gabriele Gut, and Charlotte Bunne. 2025. “AI-powered Virtual Tissues from Spatial Proteomics for Clinical Diagnostics and Biomedical Discovery.” arXiv. <https://doi.org/10.48550/arXiv.2501.06039>.
- Yuan, Zhiyuan, Wentao Pan, Xuan Zhao, Fangyuan Zhao, Zhimeng Xu, Xiu Li, Yi Zhao, Michael Q. Zhang, and Jianhua Yao. 2023. “SODB Facilitates Comprehensive Exploration of Spatial Omics Data.” *Nature Methods* 20 (3): 387–99. <https://doi.org/10.1038/s41592-023-01773-7>.