

# Mathematical Frameworks for Integrative Analysis of Multi-omics Biological Data

This manuscript ([permalink](#)) was automatically generated from [BIRSBiointegration/whitePaper@758ed6b](#) on June 26, 2020.

## Authors

---

- **Kim-Anh Lê Cao**

 [0000-0003-3923-1116](#) ·  [mixOmicsTeam](#) ·  [mixOmicsTeam](#)

Melbourne Integrative Genomics, School of Mathematics and Statistics, University of Melbourne, Australia · Funded by Grant National Health and Medical Research Council Career Development fellowship (GNT1159458)

- **Aedin C Culhane**

 [0000-0002-1395-9734](#) ·  [aedin](#) ·  [AedinCulhane](#)

Data Sciences, Dana-Farber Cancer Institute, Boston, MA, USA; Biostatistics, Harvard TH Chan School of Public Health, Boston, MA, USA · Funded by Chan Zuckerberg Initiative, NIH, DoD (need to get grant IDs)

- **Elana Fertig**

 [0000-0003-3204-342X](#) ·  [ejfertig](#) ·  [FertigLab](#)

Department of Oncology, Sidney Kimmel Comprehensive Cancer Center, Johns Hopkins University School of Medicine, Baltimore, MD, USA; Department of Biomedical Engineering, Johns Hopkins University School of Medicine, Baltimore, MD, USA; Department of Applied Mathematics and Statistics, Johns Hopkins University Whiting School of Engineering, Baltimore, MD, USA · Funded by National Institute of Health, National Cancer Institute; National Institute of Health, National Institute of Dental and Craniofacial Research; Lustgarten Foundation; Emerson Foundation; Allegheny Health Network

- **Jane Roe**

 [XXXX-XXXX-XXXX-XXXX](#) ·  [janeroe \(PLEASE COPY/PASTE DO NOT EDIT THIS ONE\)](#) ·  [XXX](#)

Department of Something, University of Whatever; Department of Whatever, University of Something · Funded by XX

# **Abstract**

---

## **Introduction to single cell and imaging multi-omics**

---

## **Current multi-omic technologies**

---

## **Challenges for interpretation**

---

**Interpretation requires a good understanding of the methods (no black boxes, a common language)**

Communicating within the field: what approaches are we talking about?

**Interpretation of the output from the analyses of the data is facilitated by the incorporation of contiguous information.**

Redundant biological knowledge and incorporation of information from databases are important in the process. Biological Interpretations: bridges to data bases such as KEGG, Gene Ontology, HCA. Validation through complementary data.

## **Visualization tools for interpretation and communication**

Examples: Brushing UMAP.

**Explaining results to biologists through generative models and simulations (ex: Factor Analysis).**

## **Issues of over-discretization, over-simplification**

Example 1: The notion of cell-type is insufficient (Communication challenge with biologists about tradeoffs between focusing on rare cell types vs. more “continuous” view on cell types).

Problem with loss of information in the desire to simplify.

## **Counterexamples**

---

## **Case studies**

---

### **scNMT-seq as a case-study for epigenetic regulation**

#### **Overview and biological question**

#### **Computational challenges**

- Identification of multi-omics signatures that characterize lineage, stage or both.

- Handling missing values
- Do epigenetic changes in some genomic contexts affect cell fate decision more than others? If so, how?

## **Methods for stats/math analyses and results summary**

### **scRNA-seq + FISH as a case study for spatial transcriptomics**

#### **Overview and biological question**

#### **Computational challenges**

- Can scRNA-seq data be overlaid onto seqFISH for resolution enhancement
- What is the minimal number of genes needed for data integration?
- Are there signatures of cellular co-localization or spatial coordinates in the non-spatial scRNA-seq data?

## **Methods for stats/math analyses and results summary**

### **Spatial proteomics and cross-study analysis**

#### **Overview and biological question**

#### **Computational challenges**

- Integrating partially-overlapping proteomic data collected on different patients with similar phenotypes
- Integration of spatial x-y coordinate co-location and co-expression
- Integration with other 'omics datasets (e.g., scRNA-seq) to support the results of these proteomic analyses
- Can we predict the spatial expression patterns of proteins measured on mass-tag but not measured in the MIBI-TOF data?
- What additional information can we learn about the different macrophage and immune populations in breast cancer by conducting integrated analyses of these datasets?

## **Methods for stats/math analyses and results summary**

### **Overview of common analytical methods spanning technologies / case studies**

---

- matrix factorization
- neural network / autoencoders

### **Software strategies to enable analyses of multimodal single cell experiments**

---

#### **Key questions**

- How should multimodal single cell data be managed for interactive and batch analyses?
- What methods will help software developers create scalable solutions for multimodal single cell analysis?

- How can we ensure that visualization methods that are central to multimodal single cell analysis are usable by researchers with visual impairments?

## Data management strategies

- Abstract data type: “multiassay experiment”. This reflects the idea that each mode will be characterized by a different collection of features on possibly non-overlapping collections of samples. The metadata on features should be clearly and conventionally defined. For example, genes and transcripts are enumerated using Ensembl catalog identifiers; regions of accessibility are defined using genomic coordinates in a clearly specified reference build. Metadata on samples must include all relevant information on experimental conditions such as treatment, protocol, and date of technical processing.

(More fodder-AS)

**Key points:** 1) What do we want to store and share? Data object vs analysis object. How would the design change based on what is stored? 2) Do we need a flexible, universal framework (e.g. MAE) or an experiment class for every possible combination of modalities or technologies? 3) Do we have adequate data representation for all “assays”?

- Multi-modal single cell data may consist of multi-assay measurements from the same cell (e.g. citeSeq, sci-CAR) or integration of multi-assay measurements from distinct cells from the same or distinct starting samples. A sample here refers to the biological specimen of origin (tissue A from individual X). A data container for a multi assay analysis must hold
  1. Assay slots containing variables or features from multiple modalities (e.g. gene expression units from sc-RNAseq and protein units from sc-proteomics). In some cases, the feature may be multidimensional (e.g. spatial coordinates, locations of eQTLs).
  2. Observations or cell identities
  3. Metadata for sample of origin for the individual cells, e.g. study, center, phenotype, perturbation.
  4. A map between the different assays to enable analysis
- The MAE is such a Bioconductor container for overlapping observations, and may serve as a starting point for further expansion. Besides the primary data elements for storing “data objects”, the summarizedExperiment class offers attributes and Methods for storing results of analysis, as an “analysis object”.
  - While common assays such as RNAseq and ATACseq have well-defined data representations (e.g. transcript names), data representation need to be defined newer assays, which may need multiple dimensions for adequate definition (e.g. x, y, z coordinates for images).
  - The observations of different modalities may not be directly comparable (e.g. RNA may be measured from individual cells but spatial transcriptomics may cover a few cells in the matched area).
  - In the absence of universal standards, the metadata may vary from analysis to analysis.

(Note, standard BioC container/class terms may not be correctly used. End of fodder - AS)

- Serializations and data access methods for
  - spatial transcriptomics
  - scNMT-seq ...

## Scalability strategies

## Reducing barriers to interpretable visualizations

[An overview of the issues with impaired color perception](#)

[US Government tools for accessibility](#)

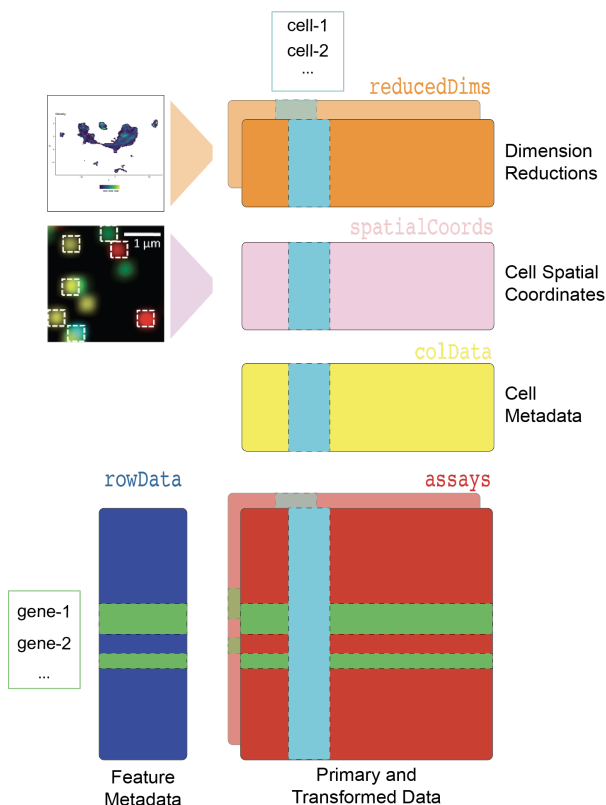
## Details of working components

Type	Brief name (link)	Description	URL	Author email
R data class	MultiAssayExperiment	unify multiple experiments	bioco nduct or.org	many
R data class	<a href="#">SpatialExperiment</a>	SpatialExperiment: a collection of S4 classes for Spatial Data	<a href="https://github.com/dri ghelli/ SpatialExper iment /">https://github.com/dri ghelli/ SpatialExper iment /</a>	...
---	-----	-----	---	-----
R package	Giotto	Spatial transcriptomics	...	...
python library	PyTorch	deep learning	...	...
R scripts	<a href="#">SpatialAnalysis</a>	Scripts for SpatialExperiment usage	<a href="https://github.com/dri ghelli/ SpatialAnaly sis">https://github.com/dri ghelli/ SpatialAnaly sis</a>	...
Matlab package	<a href="#">CytoMAP</a>	CytoMAP: A Spatial Analysis Toolbox Reveals Features of Myeloid Cell Organization in Lymphoid Tissues	<a href="https://gitlab.com/gernerlab/cytomap">https://gitlab.com/gernerlab/cytomap</a>	...
R package	cytomapper	cytomapper: Visualization of highly multiplexed imaging cytometry data in R	<a href="https://github.com/BodenmillerGroup/cytomapper">https://github.com/BodenmillerGroup/cytomapper</a>	...

Type	Brief name (link)	Description	URL	Author email
...	histoCAT	histoCAT: analysis of cell phenotypes and interactions in multiplex image cytometry data	<a href="https://github.com/BodenmillerGroup/histoCAT">https://github.com/BodenmillerGroup/histoCAT</a>	...
Python package	SpaCell	SpaCell: integrating tissue morphology and spatial gene expression to predict disease cells	<a href="https://github.com/BioMedicalMachineLearning/SpaCell">https://github.com/BioMedicalMachineLearning/SpaCell</a>	...
R package	Spaniel	Spaniel: analysis and interactive sharing of Spatial Transcriptomics data	<a href="https://github.com/RachelQueen1/Spaniel/">https://github.com/RachelQueen1/Spaniel/</a>	...
...	ST viewer	ST viewer: a tool for analysis and visualization of spatial transcriptomics datasets	<a href="https://github.com/jfnavarr/st_viewer">https://github.com/jfnavarr/st_viewer</a>	...
R package	Seurat	...	<a href="https://github.com/satijalab/seurat">https://github.com/satijalab/seurat</a>	...
Python package	Scanpy	<a href="https://github.com/theislab/scanpy">https://github.com/theislab/scanpy</a>	...	
R package	SpatialLIBD	Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex	<a href="https://github.com/LieberInstitute/spatialLIBD">https://github.com/LieberInstitute/spatialLIBD</a>	...
Shiny app	Dynverse	A comparison of single-cell trajectory inference methods: towards more accurate and robust tools	<a href="https://zouter.shinyapps.io/server/">https://zouter.shinyapps.io/server/</a>	...

Type	Brief name (link)	Description	URL	Author email
R package	Cardinal	Cardinal: an R package for statistical analysis of mass spectrometry-based imaging experiments	<a href="https://cardinal.msi.org/">https://cardinal.msi.org/</a>	...
R package	CoGAPS	<a href="https://doi.org/10.1016/j.cels.2019.04.004">https://doi.org/10.1016/j.cels.2019.04.004</a>	<a href="https://github.com/FertigLab/CoGAPS">https://github.com/FertigLab/CoGAPS</a>	...
R package	projectR	<a href="https://doi.org/10.1016/j.cels.2019.04.004">https://doi.org/10.1016/j.cels.2019.04.004</a>	<a href="https://github.com/gensofeve/projectR">https://github.com/gensofeve/projectR</a>	...
---	-----	-----	---	-----
ExperimentHub package	SingleCellMultiModal	Serves multiple datasets obtained from GEO and other sources and represents them as MultiAssayExperiment objects	<a href="https://biocductor.org">biocductor.org</a>	...
Proteomics&Transcriptomics data	hubmap	...	<a href="https://hubmapconsortium.org/">https://hubmapconsortium.org/</a>	...
Protein atlas	Protein Atlas	...	<a href="https://www.proteinatlas.org">https://www.proteinatlas.org</a>	...
Proteome database	The proteome landscape of the kingdoms of life	...	<a href="http://www.proteomesoflife.org">www.proteomesoflife.org</a>	...
Mass Spectrometry archive	MassIVE.quant	...	...	<a href="https://massive.ucsd.edu/ProteoSAFe/static/massive-quant.jsp">https://massive.ucsd.edu/ProteoSAFe/static/massive-quant.jsp</a>

Here is the schematic of SpatialExperiment class from Dario Righelli.



## Techniques and challenges for benchmarking methods

Outline:

- We must first define what we are benchmarking
  - recovery of cell types / clusters
  - discovery of relationships between data modalities, e.g. gene regulatory relationships observed between chromatin accessibility and gene expression
  - ...
- Strategies for benchmarking
  - simulation (and we can discuss the difficulties with simulating covariance structure across features and data modalities)
  - benchmarking datasets
  - cross-validation within study (and we can discuss issues in matching dimensions of latent space across folds). For this Mike has a lot of literature in Google Doc to include on papers that have performed either permutation or cross-validation to assess model performance.
  - cross-study validation (are relationships discovered in one dataset present in other datasets, potentially looking across single cell and bulk)

## Discussion

### Emerging analytical methods and technologies

### Community needs for data structures, analysis methods, etc

## Glossary

Consensus term	Synonyms	Description
----------------	----------	-------------



Consensus term	Synonyms	Description
Network	Graph	A set of <i>nodes</i> , representing objects of interest, linked by <i>edges</i> , representing specific relationships between nodes.
Node	Vertex	Element of interest in a network and linked to other nodes. For example: people, cells, proteins or genes. Nodes can have several properties called <i>attributes</i> like cell type or position.
Edge	Link	The relationship between 2 nodes in a network. For example: friendship in social networks, cells in contact in a spatial network, or gene-gene interactions in a gene regulatory network.

## References

---