

Community-wide hackathons establish foundations for emerging single cell data integration

This manuscript ([permalink](#)) was automatically generated from [BIRSBiointegration/whitePaper@1976994](#) on July 2, 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

Comprehensive characterization of biological systems with multi-omics

- Single cell community has advanced technologies to enable concurrent processing of biological systems at multiple molecular resolutions
- The lack of prior knowledge and gold standard benchmark naturally leads to a data-driven approach

New single cell multi omics initiatives:

- Human Cell Atlas (HCA): assess variation in normal tissues
- Brain initiative and Allen Brain
- Human Tumor Atlas Network (HTAN): Single-cell, longitudinal, and clinical outcomes atlases of cancer transitions for diverse tumor types.

What bulk multi-omics (e.g. TCGA, ENCODE) have taught us:

- Type of omics that can answer a specific biological question
- The value of open resources for methodological developments
- New hypotheses

Using hackathons to illustrate analysis standards and challenges for capturing biological information from multi-omics technologies

- Brief overview of our three hackathon studies highlighting state of the art challenges (e.g., spatial transcriptomics, cross-study analysis, epigenetic regulation)
- Challenges include issues with noise and experimental design, Time lag between regulatory levels not addressed and many open questions remain (e.g methylation / gene expression), Direction of regulation not captured
- We present our findings from hackathon case studies that helped us obtain benchmarks and define a common language for multi-omics
- **Objectives of this paper**
 - Provide guidelines on tools / data / technologies / methods and needs to model the multi-scale regulatory processes in biological systems for a computational biologist audience
- **Outline and messages**
 - Cellular and molecular regulation is fundamentally multi-scale and captured by distinct data modalities
 - Traditional hypothesis-driven multi-omics/view studies only consider one facet of these technologies, but more can be learned through a holistic approach extending into atlases
 - We present our findings from hackathon case studies that helped us obtain a broader picture and language

Outline of the paper:

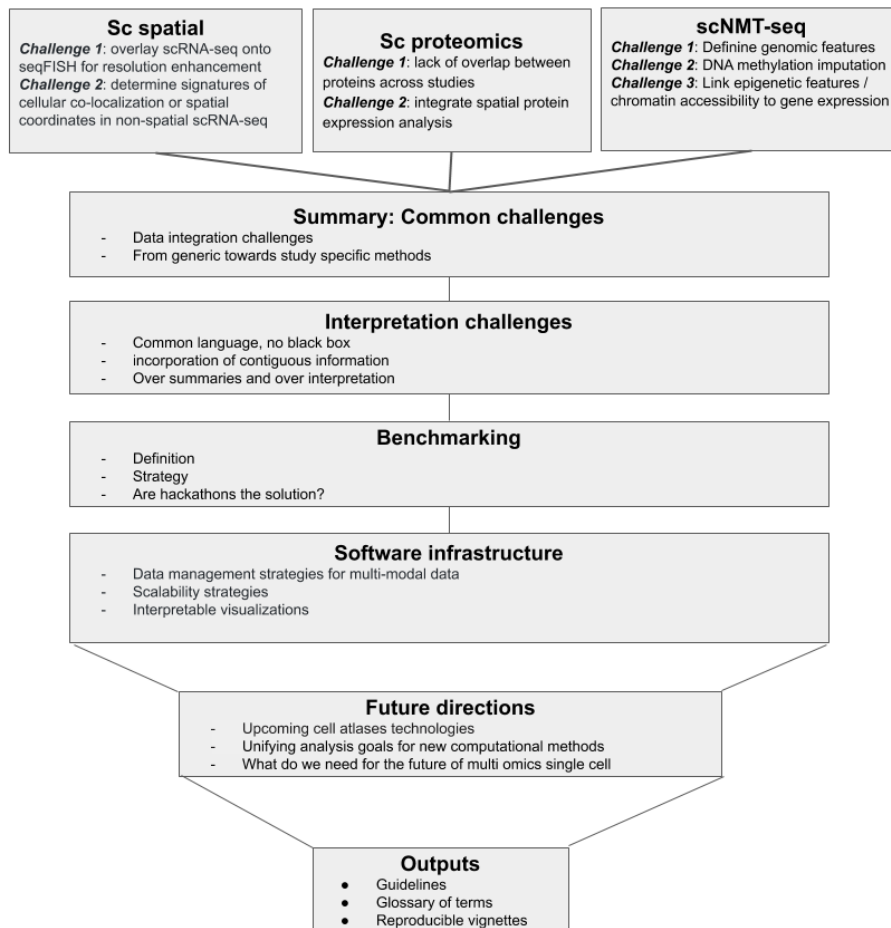


Figure caption: Main challenges discussed during our brainstorming sessions from the hackathons

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

Overview and biological question

Computational challenges

Challenge 1: overlay of scRNA-seq onto seqFISH for resolution enhancement

Challenge 2: determine signatures of cellular co-localization or spatial coordinates in non-spatial scRNA-seq

Spatial proteomics and cross-study analysis

Overview and biological question

Computational challenges

Challenge 1: address the lack of overlap between proteins across studies

Challenge 2: spatial protein expression analysis

scNMT-seq as a case-study for epigenetic regulation

Overview and biological question

Computational challenges

Challenge 1: defining genomic features

Challenge 2: DNA methylation imputation

Challenge 3: Linking epigenetic features / chromatin accessibility to gene expression

Analytical approaches for hackathons and commonalities for multi-omics analysis illustrated by the series of hackathons [Figure 5 + Table 1]

Short introduction explaining that we focus here on the common challenges across hackathons

Summary of hackathon study-specific methods

- Table describes method, foundation in the context of previous bulk and single cell literature, and technology dependence
 - Attempts to tweak existing methods and challenges associated in hackathons
 - List methods that are either technology dependent (e.g. spatial) vs universal and how to choose them

Dependence on pre-processing method and/or variable selection

- These steps are key and affect downstream analyses
 - Normalization / data transformation (seqFish), pre-processing, gene summaries (scNMT-seq) to variable selection (seqFish)
 - reproducibility difficult / no consensus. e.g. Alexis selected 19 genes whereas Zhu original paper based on 47 genes (difference in methods / processed data)
- Hackathon data pre-processed enable better comparisons across methods
- No consensus reached as those are emerging data with no ground truth nor established biological results

Approaches for partial overlap of information (cells / features) and how to predict (cell type, dataset) using another data set

- Overlap in each study
 - seqFish: same features but not cells; scProt: same proteins, not cells but similar patients; scNMT-seq: same cells but not features
 - How it was solved (Table)
- Anchoring information across datasets or studies is needed (Figure)
- Incorporation of existing biological knowledge

- 'From discovery to detection' (Meuleman + debrief), time is ripe to include more knowledge in our data driven approaches
- Challenge: Partial cell overlap (but no features matching) and No overlap were not addressed

Managing differences in scale and size for datasets that do not match cells or features

- Hackathons datasets did not match cells or features.
 - scNMTQ-seq: MOFA limitation when # features vary (and size of datasets).
 - seqFish: greedy approach to select the best gene subset (Alexis, size); consider batch effect removal method (Amrit, scale)
- Consensus on projection based methods, even if pre-processing was applied (Table)
- Additional weighting is needed (e.g. Arora, Abadi).

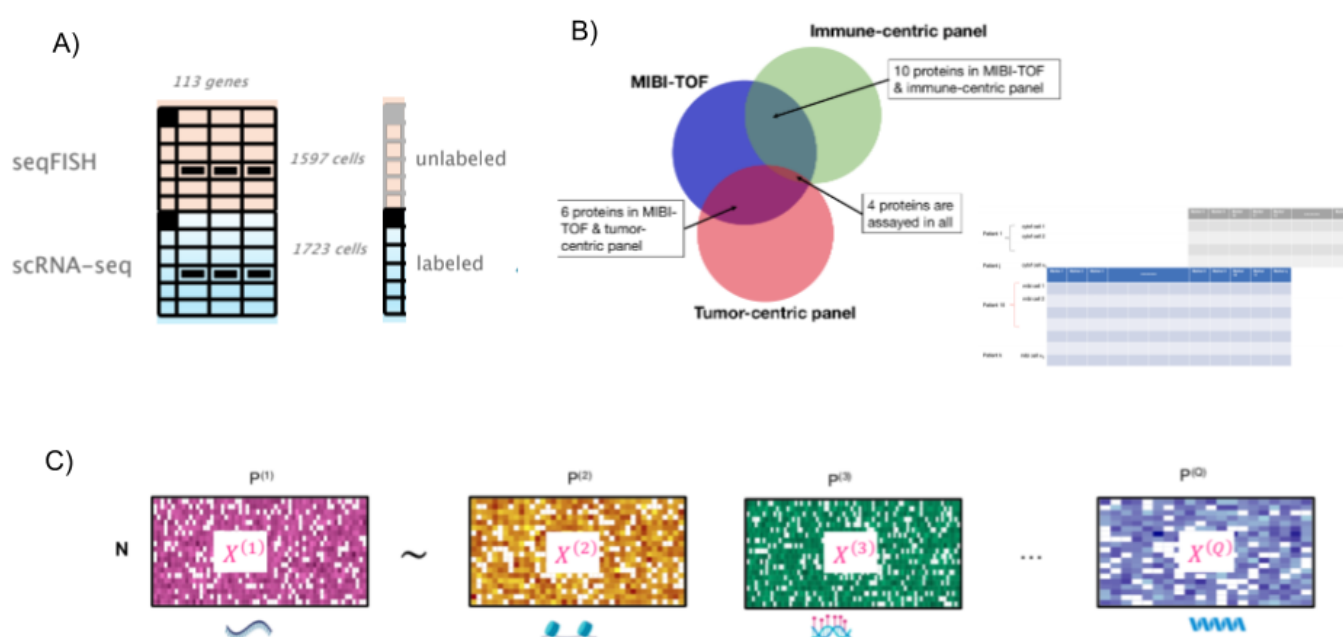


Figure caption:

- scSpatial: required overlap of features (genes), but cells do not overlap. Cell type prediction for seqFISH data was performed based on scRNA-seq (known) [credit: Amrit Singh]
- scProteomics: a small number of features overlap (proteins) but patients and cells do not overlap. Data imputation (?), spatial co-localization or cell type prediction was performed [credit: Lauren Hsu and Pratheepan Jenagan]
- scNMT-seq: cells are matching across assays but features do not overlap. Data integration was performed [credit: Al Abadi]

Challenges for interpretation

Interpretation for quantitative team: requires a good understanding of the methods (no black boxes, a common language)

Communicating within the field: what approaches are we talking about, this requires agreement on a glossary (ref: Table1). Supervised versus unsupervised methods. Visualization figures are useful for the mathematical and data science team in the Explorator phase.

Interpretation for biologists:

Understanding 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 workflow. Biological interpretations are facilitated by bridges to databases such as KEGG, Gene Ontology, Human Cell Atlas, Biomart and many other databases. Validation through complementary data and sequential experimental design.

Visualization tools for interpretation and communication to biologists

There are pitfalls in using sophisticated graphics which lead to over-interpretation or misinterpretation (size of clusters in tSNE related to sampling baselines rather than density, ...) Example of effective visual interpretation tools : brushing UMAP (Kris Sankaran).

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

Several difficulties arise when explaining summaries and conclusions, problems encountered include non-identifiability of models or non-sufficiency of summaries, simulations can often provide effective communication tools.

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

Techniques and challenges for benchmarking methods

We must first define what we are benchmarking

- Often the goal in benchmarking is recovery of known cell types with processing of raw data, quantification, and clustering. The Adjusted Rand Index (ARI) or other metrics for partitions are used.
- One may also attempt to benchmark methods for their ability to discover known relationships between data modalities, e.g. gene regulatory relationships observed between chromatin accessibility and gene expression. However, this is made difficult by the fact that these relationships are not fully known at the single cell level.

Strategies for benchmarking

- Simulation is useful for having known truth, but it is difficult to simulate realistic covariance structure across features and across data modalities.
- Benchmarking datasets (add examples from Google Doc). Benchmark datasets for single cell studies have largely centered around measuring sequencing depth and diversity of cell types derived from a single assay of interest (e.g. scRNAseq). A benchmark dataset serves a few purposes:
 - Provides ground truth for the intended effect of exposure in a proposed study design.

- Provides validation for a data integration task for which a new computational method may be proposed.

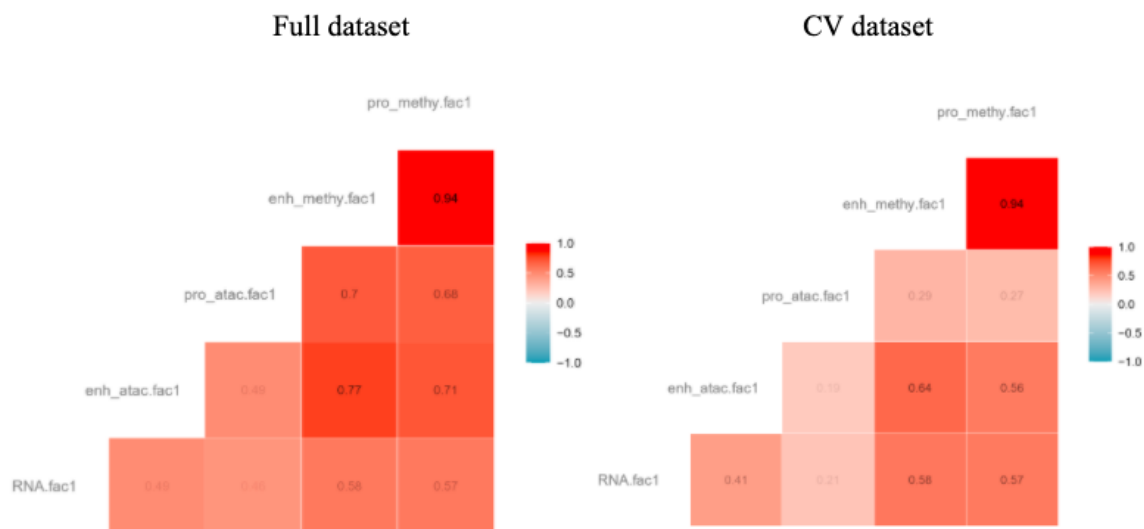
For multi-modal assays, while the intended effects can vary based on the leading biological questions, one may abstract out common data integration tasks such as co-embedding, mapping or correlation, and inferring causal relationships. We distinguish data integration from further downstream analyses that may occur on integrated samples such as differential analysis of both assays with regard to a certain exposure.

Both the intended effects and data integration task rely on study design that takes into account:

- Biological and technical variability via replicates, block design, and randomization.
- Power analysis for the intended effect or data integration task.
- Dependencies between modalities, for e.g. gene expression depending on gene regulatory element activity, requires that experiment design must also account for spatial and temporal elements in sampling for a given observation.

As such, no universal benchmark data scheme may suit every combination of modality, and benchmark datasets may be established for commonly used combinations of modalities or technologies, towards specific data integration tasks.

- Cross-validation within study can be performed. For example the following cross-validation analysis of the scNMT-seq dataset was performed using MOFA+



A challenge with within study cross-validation is how to match dimensions of latent space across folds. (add examples from Google Doc of papers that have performed either permutation or cross-validation to assess model performance)

- Cross-study validation would assess if relationships discovered in one dataset present in other datasets, potentially looking across single cell and bulk.

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. CITE-seq, 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 scRNA-seq 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 RNA-seq and ATAC-seq 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.
 - It is crucial that data containers use consistent assay access methods (possibly through methods inheritance. e.g. from `SummarizedExperiment`). This will ensure less redundancy in development process and allow powerful implementation strategies.

(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

Color is a powerful data visualization tool that helps representing the different dimensions of our increasingly complex and rich scientific data. Color vision deficiencies affect a substantial portion of the population. Therefore, it is desirable to aim towards presenting scientific information in a manner that is as accessible as possible for all readers. Color vision deficiency leads to difficulties in perceiving patterns (the basis for the Ishihara's color vision tests) in multi-colored figures. In rare cases, the perceived patterns; e.g. in heatmaps and reduced dimension plots, can differ between individuals with normal and color deficient vision.

One strategy to address these issues is to include colorblind friendly visualizations [1] as a default setting in our visualizations. Several colorblind-friendly palettes exist (e.g., see R packages [viridis](#) and [dittoSeq](#)) and can be integrated into data presentation as the default option. Even with these palettes in place, it is desirable to limit the number (about 8-10 at a maximum) of colors in visualizations. To reduce the dependence on colors, one solution would be to include additional visual cues to differentiate regions (hatched areas) or cells (point shapes). Overall, a broader discussion regarding the accessibility of our figures that is not just limited to color vision deficiencies would be greatly beneficial towards improving data accessibility. Perhaps one tool to address broader accessibility could be the inclusion an "accessibility caption" accompanying figures which "guide" the reader's perception of the images.

[Reference 1: Color Coding](#)

[Reference 2: Points of View: Color Blindness](#)

[Viridis Color Palettes](#)

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

[US Government tools for accessibility](#)

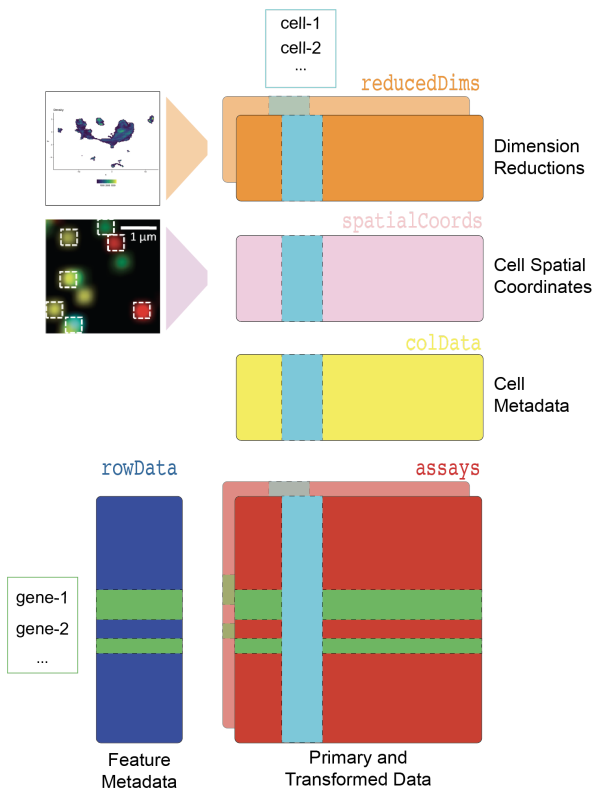
Details of working components – trimmed

you can interact with underlying data at [google sheet](#)

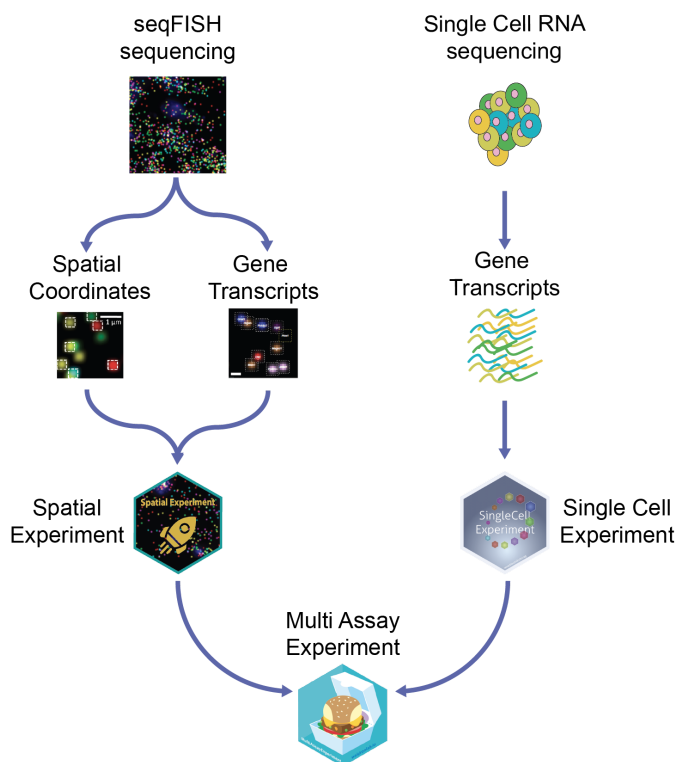
Type	Brief name (link)	Description
Matlab package	CytoMAP	CytoMAP: A Spatial Analysis Toolbox Reveals Features of Myeloid Cell Organization in Lymphoid Tissues
Matlab package	histoCAT	histoCAT: analysis of cell phenotypes and interactions in multiplex image cytometry data
Python library	PyTorch	General framework for deep learning
Python package	SpaCell	SpaCell: integrating tissue morphology and spatial gene expression to predict disease cells
Python package	Scanpy	Python package for single cell analysis
R data class	MultiAssayExperiment	unify multiple experiments

Type	Brief name (link)	Description
R data class	SpatialExperiment	SpatialExperiment: a collection of S4 classes for Spatial Data
R package	Giotto	Spatial transcriptomics
R package	cytomapper	cytomapper: Visualization of highly multiplexed imaging cytometry data in R
R package	Spaniel	Spaniel: analysis and interactive sharing of Spatial Transcriptomics data
R package	Seurat	R toolkit for single cell genomics
R package	SpatialLIBD	Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex
R package	Cardinal	Cardinal: an R package for statistical analysis of mass spectrometry-based imaging experiments
R package	CoGAPS	scCoGAPS learns biologically meaningful latent spaces from sparse scRNA-Seq data
R package	projectR	ProjectR is a transfer learning framework to rapidly explore latent spaces across independent datasets
R package	SingleCellMultiModal	Serves multiple datasets obtained from GEO and other sources and represents them as MultiAssayExperiment objects
R scripts	SpatialAnalysis	Scripts for SpatialExperiment usage
Self-contained GUI	ST viewer	ST viewer: a tool for analysis and visualization of spatial transcriptomics datasets
Shiny app	Dynverse	A comparison of single-cell trajectory inference methods: towards more accurate and robust tools
R package	mixOmics	R toolkit for multivariate analysis of multi-modal data

Here is the schematic of SpatialExperiment class from Dario Righelli.



Here is the schematic of how seqFISH data are stored in the SingleCellMultiModal package from Dario Righelli.



Discussion

Emerging analytical methods and technologies

Community needs for data structures, analysis methods, etc

Glossary

Table 1: Glossary of interchangeable terms in the field of single-cell and bulk multi-omics (multi-source) data analysis.

Consensus Term	Related Terms	Description	Citation
network	graph, adjacency matrix	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.	?
concordant	concordant, coherent, consistent	?	2
contributions	variable weights, loadings, eigenvector, axis, direction, dimension, coefficients, slopes	Contributions of the original variables in constructing the components.	3 , 4
latent factors	variates, scores, projections, components, latent/hidden/unobserved variables/factors	Weighted linear combinations of the original variables.	3 , 4
multimodal	Multiview, multiway arrays, multimodal, multidomain, multiblock, multitable, multi-omics, multi-source data analysis methods, N-integration	Methods pertaining to the analysis of multiple data matrices for the same set of observations.	3 , 5 , 6
conjoint analysis	conjoint analysis, P-integration, meta-analysis, multigroup data analysis	Methods pertaining to the analysis of multiple data matrices for the same set of variables.	3 , 4 , 7
variable	feature, variable	A measurable quantity that describes an observation's attributes. Variables from different modalities include age, sex, gene or protein abundance, single nucleotide variants, operational taxonomic units, pixel intensity <i>etc.</i>	?

Consensus Term	Related Terms	Description	Citation
biomarker	marker, biomarker	A variable that is associated with normal or disease processes, or responses to exposures, or interventions. Any change in this variable is also associated with a change in the associated clinical outcome. These variables may be used for diagnostic, monitoring, Pharmacodynamic responses. Examples include LDL cholesterol, CD4 counts, hemoglobin A1C.	8
panel	biomarker panel, biomarker signature	A subset of the originally measured variables that are determined to be associated with the outcome or response variable. This may be determined using statistical inference, feature selection methods, or machine/statistical learning.	9 , 10
observation	sample, observation, array	A single entity belonging to a larger grouping. Examples include patients, subjects, participants, cells, biological sample, usually the unit of observation on which the variables are measured <i>etc.</i>	?

References

1. Points of view: Color blindness

Bang Wong

Nature Methods (2011-06-01) <https://www.nature.com/articles/nmeth.1618>

DOI: [10.1038/nmeth.1618](https://doi.org/10.1038/nmeth.1618)

2. Consistency and overfitting of multi-omics methods on experimental data

Sean D McCabe, Dan-Yu Lin, Michael I Love

Briefings in Bioinformatics (2019-07-04) <https://doi.org/gghpmf>

DOI: [10.1093/bib/bbz070](https://doi.org/10.1093/bib/bbz070) · PMID: [31281919](https://pubmed.ncbi.nlm.nih.gov/31281919/)

3. mixOmics: An R package for 'omics feature selection and multiple data integration

Florian Rohart, Benoît Gautier, Amrit Singh, Kim-Anh Lê Cao

PLOS Computational Biology (2017-11-03) <https://doi.org/gcj84s>

DOI: [10.1371/journal.pcbi.1005752](https://doi.org/10.1371/journal.pcbi.1005752) · PMID: [29099853](https://pubmed.ncbi.nlm.nih.gov/29099853/) · PMCID: [PMC5687754](https://pubmed.ncbi.nlm.nih.gov/PMC5687754/)

4. Multivariate analysis of multiblock and multigroup data

A. Eslami, E. M. Qannari, A. Kohler, S. Bougeard

Chemometrics and Intelligent Laboratory Systems (2014-04) <https://doi.org/f52wrr>

DOI: [10.1016/j.chemolab.2014.01.016](https://doi.org/10.1016/j.chemolab.2014.01.016)

5. Multitable Methods for Microbiome Data Integration

Kris Sankaran, Susan P. Holmes

Frontiers in Genetics (2019-08-28) <https://doi.org/gf8dqn>

DOI: [10.3389/fgene.2019.00627](https://doi.org/10.3389/fgene.2019.00627) · PMID: [31555316](https://pubmed.ncbi.nlm.nih.gov/31555316/) · PMCID: [PMC6724662](https://pubmed.ncbi.nlm.nih.gov/PMC6724662/)

6. Dimension reduction techniques for the integrative analysis of multi-omics data

Chen Meng, Oana A. Zeleznik, Gerhard G. Thallinger, Bernhard Kuster, Amin M. Gholami, Aedín C. Culhane

Briefings in Bioinformatics (2016-07) <https://doi.org/f83qvd>

DOI: [10.1093/bib/bbv108](https://doi.org/10.1093/bib/bbv108) · PMID: [26969681](https://pubmed.ncbi.nlm.nih.gov/26969681/) · PMCID: [PMC4945831](https://pubmed.ncbi.nlm.nih.gov/PMC4945831/)

7. Robust meta-analysis of gene expression using the elastic net

Jacob J. Hughey, Atul J. Butte

Nucleic Acids Research (2015-07-13) <https://doi.org/f7nnbm>

DOI: [10.1093/nar/gkv229](https://doi.org/10.1093/nar/gkv229) · PMID: [25829177](https://pubmed.ncbi.nlm.nih.gov/25829177/) · PMCID: [PMC4499117](https://pubmed.ncbi.nlm.nih.gov/PMC4499117/)

8. Biomarker definitions and their applications

Robert M Califf

Experimental Biology and Medicine (2018-02-06) <https://doi.org/gcxh8n>

DOI: [10.1177/1535370217750088](https://doi.org/10.1177/1535370217750088) · PMID: [29405771](https://pubmed.ncbi.nlm.nih.gov/29405771/) · PMCID: [PMC5813875](https://pubmed.ncbi.nlm.nih.gov/PMC5813875/)

9. Biomarker signatures of aging

Paola Sebastiani, Bharat Thyagarajan, Fangui Sun, Nicole Schupf, Anne B. Newman, Monty Montano, Thomas T. Perls

Aging Cell (2017-04) <https://doi.org/d2cm>

DOI: [10.1111/accel.12557](https://doi.org/10.1111/accel.12557) · PMID: [28058805](https://pubmed.ncbi.nlm.nih.gov/28058805/) · PMCID: [PMC5334528](https://pubmed.ncbi.nlm.nih.gov/PMC5334528/)

10. Biomarker Panels in Critical Care

Susan R. Conway, Hector R. Wong

Critical Care Clinics (2020-01) <https://doi.org/d2cn>
DOI: [10.1016/j.ccc.2019.08.007](https://doi.org/10.1016/j.ccc.2019.08.007) · PMID: [31733684](https://pubmed.ncbi.nlm.nih.gov/31733684/)