



Advanced topics

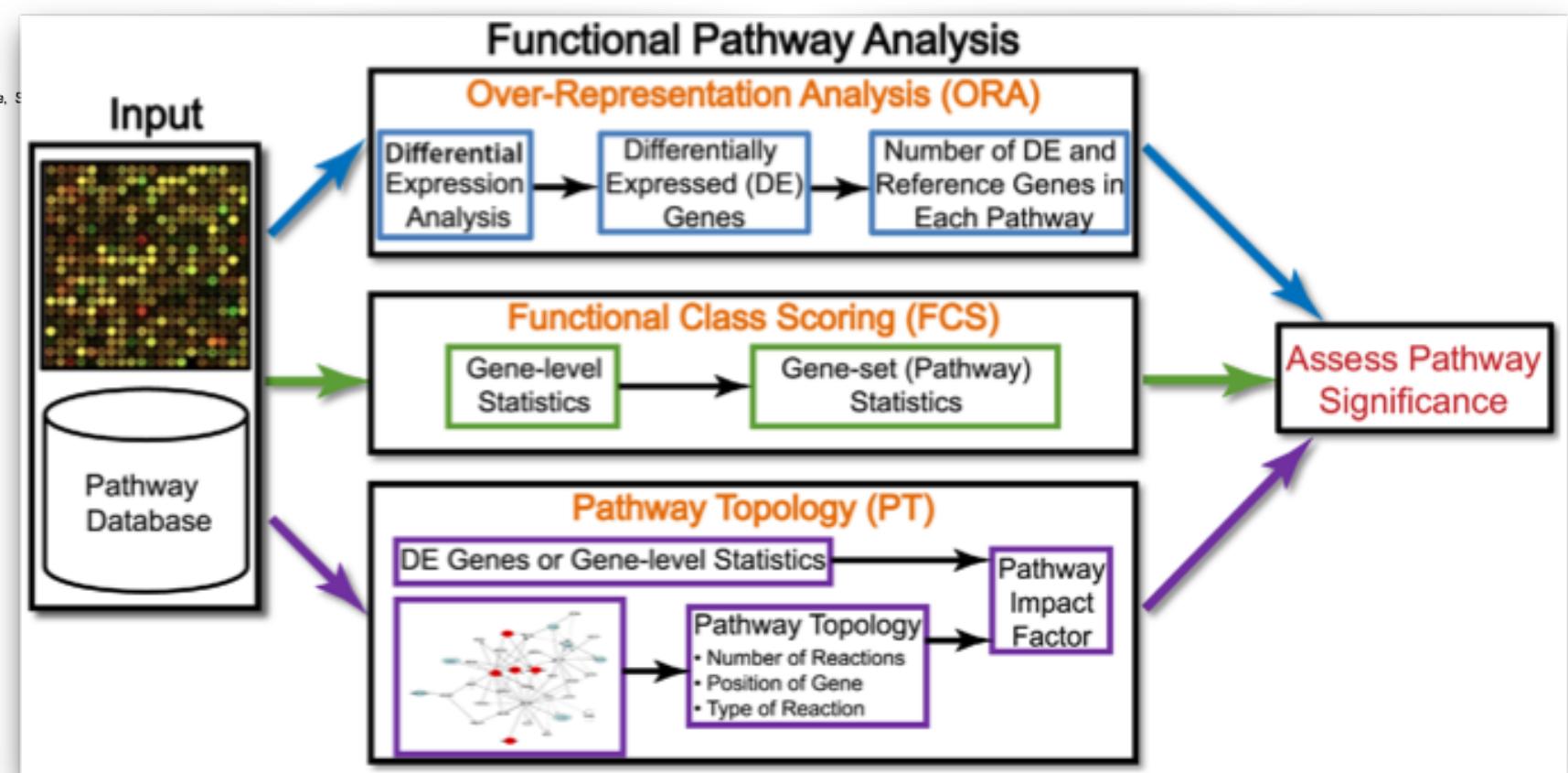
- Geneset testing
- Single cell RNA-seq



Ten Years of Pathway Analysis: Current Approaches and Outstanding Challenges

Purvesh Khatri^{1,2*}, Marina Sirota^{1,2}, Atul J. Butte^{1,2*}

¹ Division of Systems Medicine, Department of Pediatrics, Stanford University School of Medicine, S Children's Hospital, Palo Alto, California, United States of America





Casting differential expression onto biological knowledge: Functional category analysis versus gene set analysis

Motivation: DE genes might belong to a known pathway or might be the top genes from a related experiment; gene set as a whole might be altered, even if individual genes are not.

Starting point:	threshod, set of DE genes	gene-level statistics
Tool examples:	DAVID [C] goseq [C]	GSEA [S] roast [S] CAMERA [C]

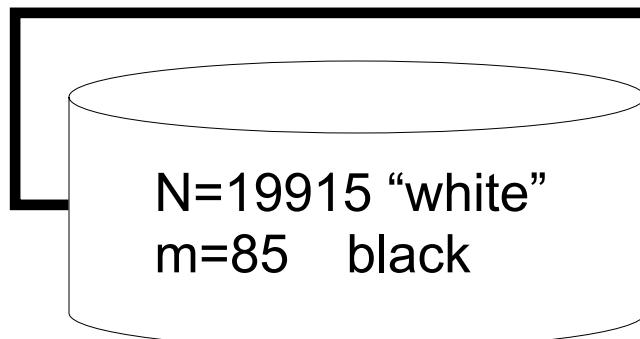
S = self-contained
C = competitive



Functional category analysis: Overlap statistics

Question: Say you have a set of 85 genes (of a total 20000 genes) known to be associated with some function. Calculate the probability of randomly selecting 40 or more (overrepresented) of those genes in a list of 100 DE genes.

Answer: Hypergeometric (i.e. the “urn” problem).



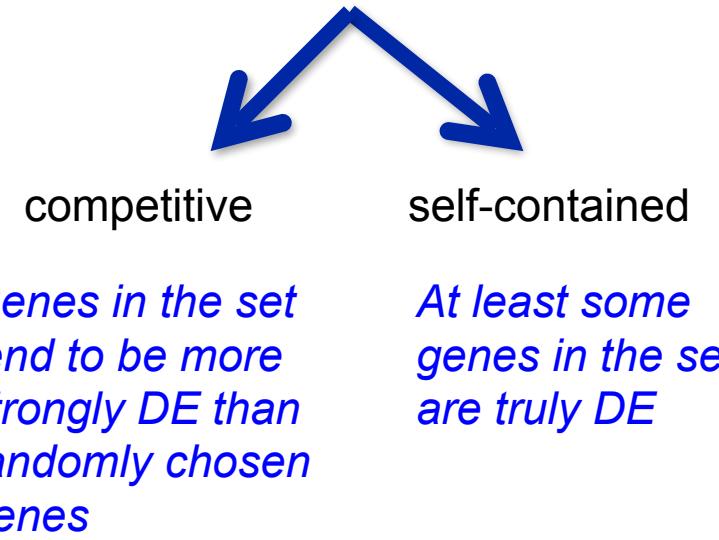
↓
n=100
k=40

$$P(X = k) = \frac{\binom{m}{k} \binom{N-m}{n-k}}{\binom{N}{n}}.$$

e.g. FunSpec (yeast) - Robinson et al. 2002 BMC BioRxiv; DAVID; topGO



Gene set analysis: what is the hypothesis (test)?



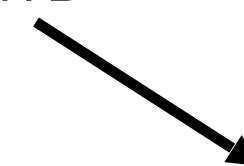


Viewing gene sets

Cell adhesion genes



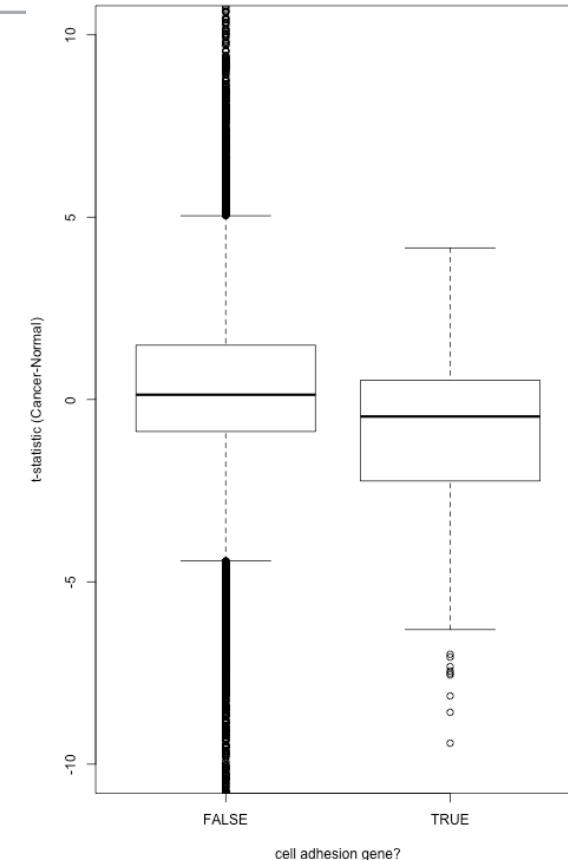
Genes regulated by MYB



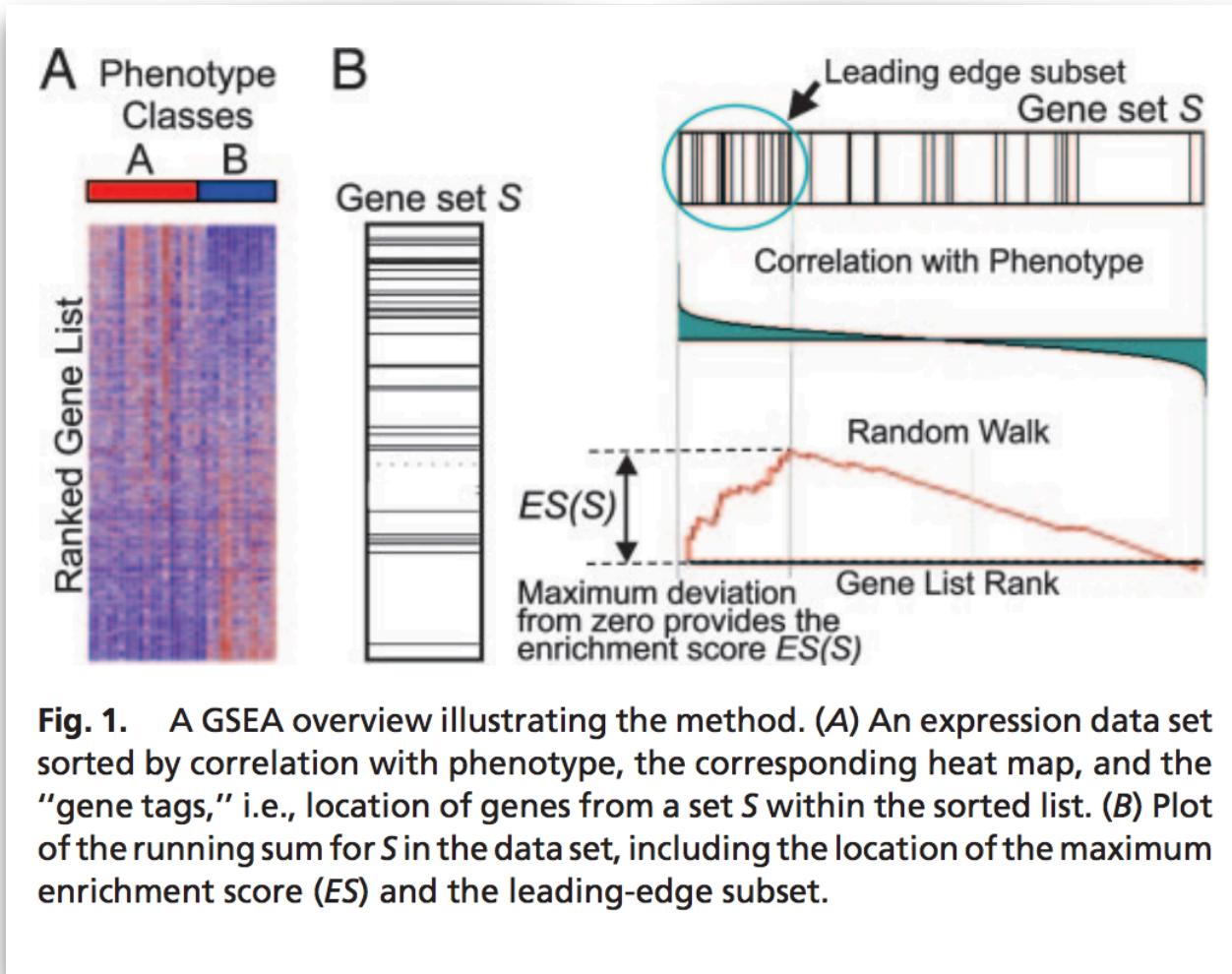
Positive



Negative



Gene set enrichment analysis (GSEA)



Self-contained.

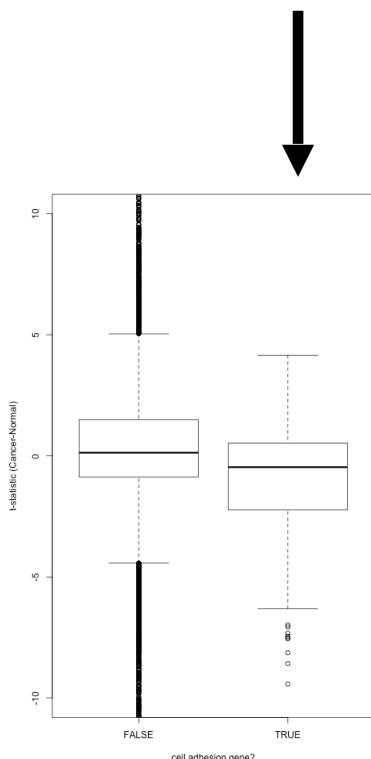
Permutation P-value:
Sample permutation is
done, which preserves
gene correlation.

But, it has limited use in
small samples (i.e. very
few possible
permutations).

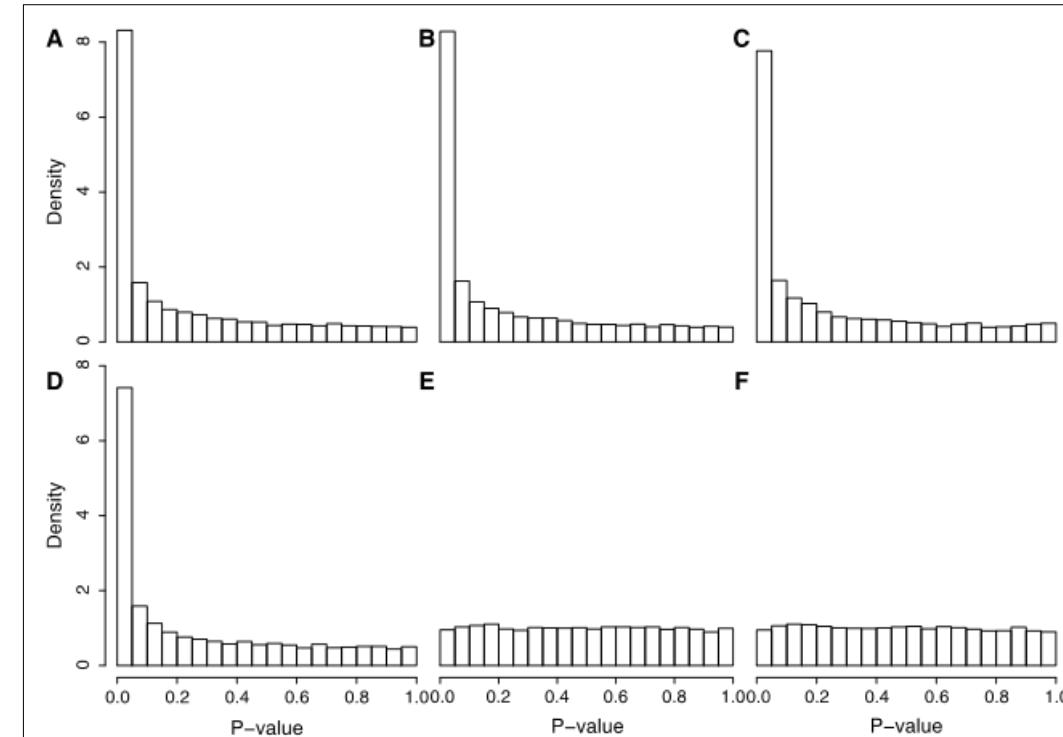
Now switches to a gene-based permutation (competitive) in small samples.

CAMERA (Correlation Adjusted MEan RAnk)

Cell adhesion genes



Main criticism of (naïve, gene-permutation) **competitive** tests is that the correlation structure is broken.



Distributions of p-value:
no differential expression

A geneSetTest
B geneSetTest [r]
C sigPathway
D PAGE
E CAMERA
F CAMERA [r]



How much of this is storytelling?

A Critical Assessment of Storytelling: Gene Ontology Categories and the Importance of Validating Genomic Scans

Pavlos Pavlidis,^{*1} Jeffrey D. Jensen,² Wolfgang Stephan,³ and Alexandros Stamatakis¹

¹The Exelixis Lab, Scientific Computing Group, Heidelberg Institute for Theoretical Studies (HITS gGmbH), Heidelberg, Germany

²Ecole Polytechnique Fédérale de Lausanne, School of Life Sciences, Lausanne, Switzerland

³Section of Evolutionary Biology, Biocenter, University of Munich, Planegg-Martinsried, Germany

***Corresponding author:** E-mail: pavlidisp@gmail.com.

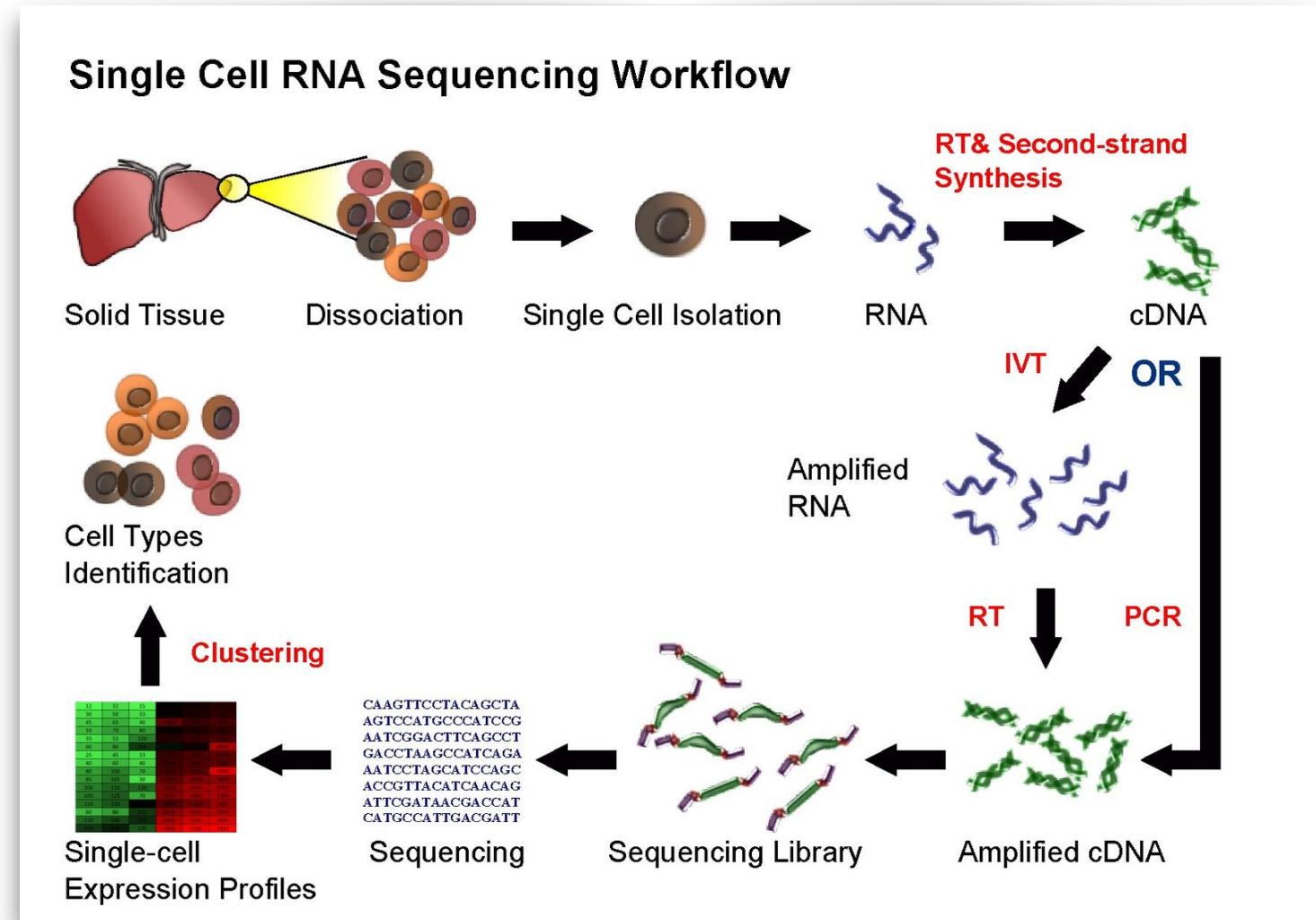
Associate editor: Arndt von Haeseler

Abstract

In the age of whole-genome population genetics, so-called genomic scan studies often conclude with a long list of putatively selected loci. These lists are then further scrutinized to annotate these regions by gene function, corresponding biological processes, expression levels, or gene networks. Such annotations are often used to assess and/or verify the validity of the genome scan and the statistical methods that have been used to perform the analyses. Furthermore, these results are frequently considered to validate “true-positives” if the identified regions make biological sense *a posteriori*. Here, we show that this approach can be potentially misleading. By simulating neutral evolutionary histories, we demonstrate that it is possible not only to obtain an extremely high false-positive rate but also to make biological sense out of the false-positives and construct a sensible biological narrative. Results are compared with a recent polymorphism data set from *Drosophila melanogaster*.

Key words: genome scanning, positive selection, gene ontology, validation, literature mining.

https://en.wikipedia.org/wiki/Single_cell_sequencing





Raphael Gottardo

@raphg

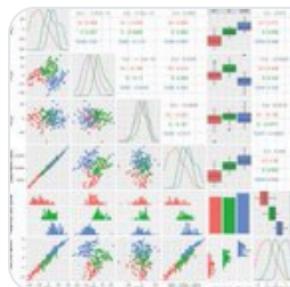


Following



A reminder that experimental design is critical for single-cell analysis, specifically scRNA-seq. Engage with a statistician as early as possible!

nature.com/articles/nbt.3120 You don't want #confounding between technical and biological effects @humancellatlas @Bioconductor

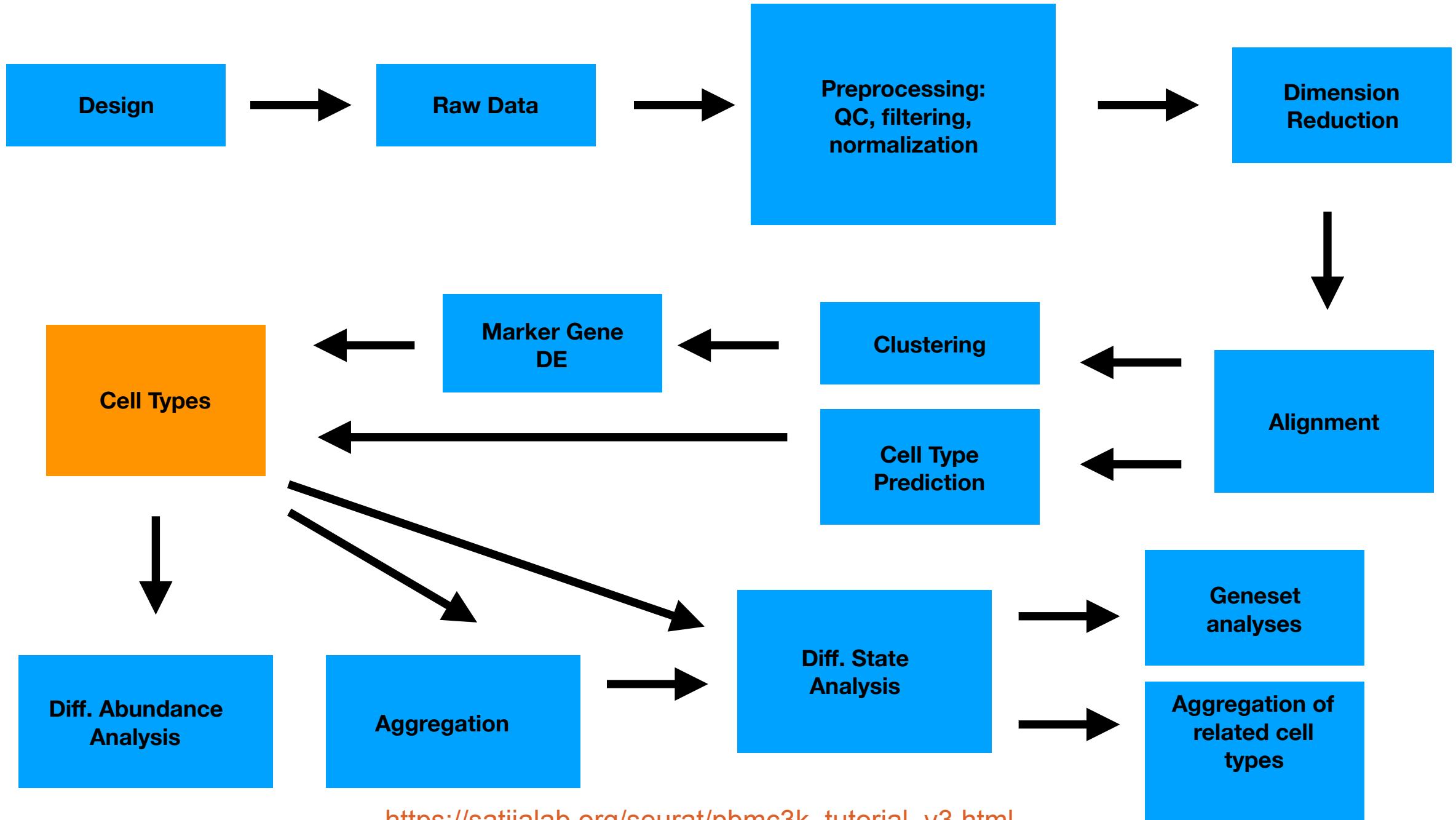


The contribution of cell cycle to heterogeneity in single-cell...

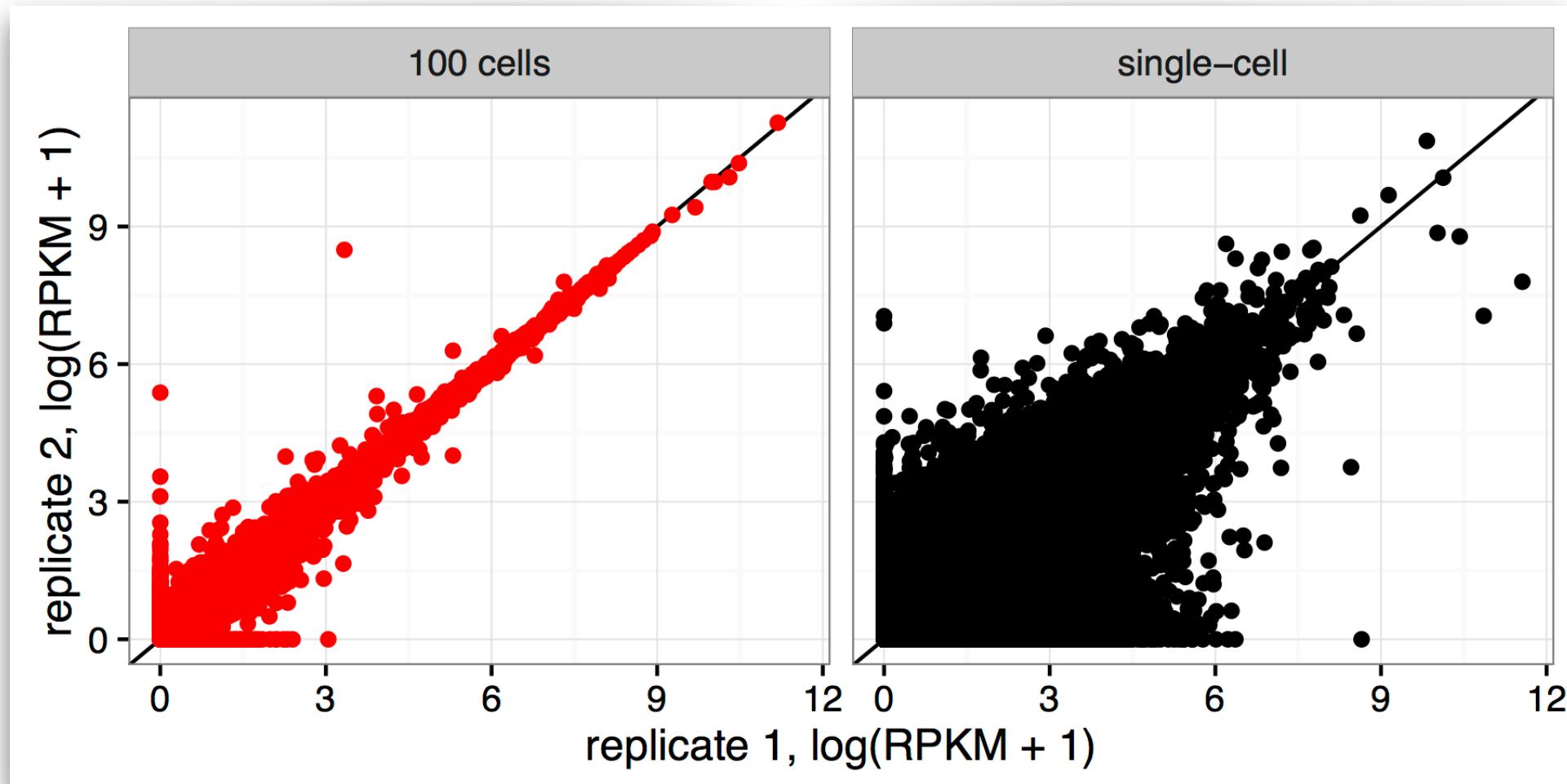
The contribution of cell cycle to heterogeneity in single-cell RNA-seq data

nature.com

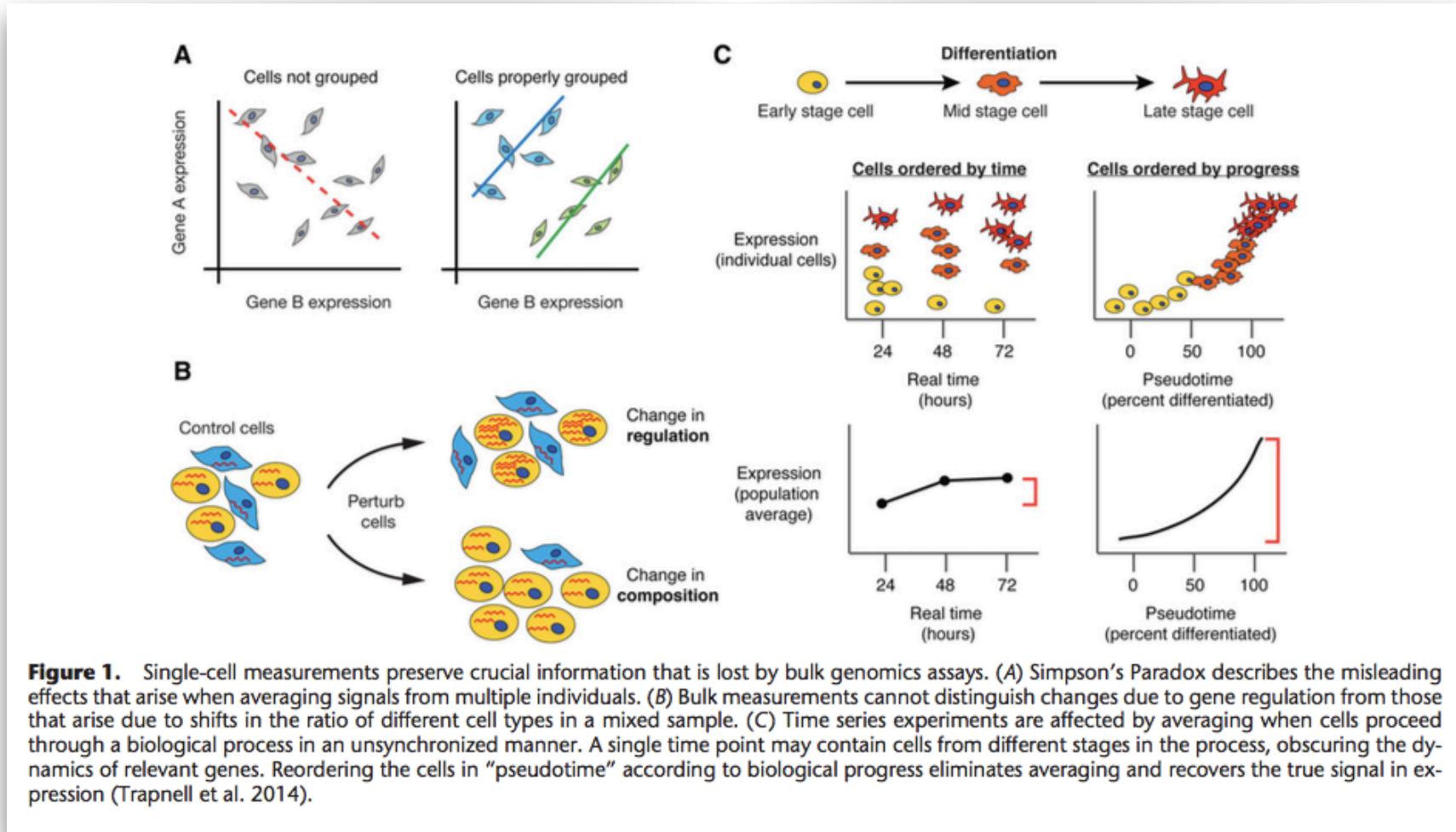
11:00 AM - 11 Dec 2018



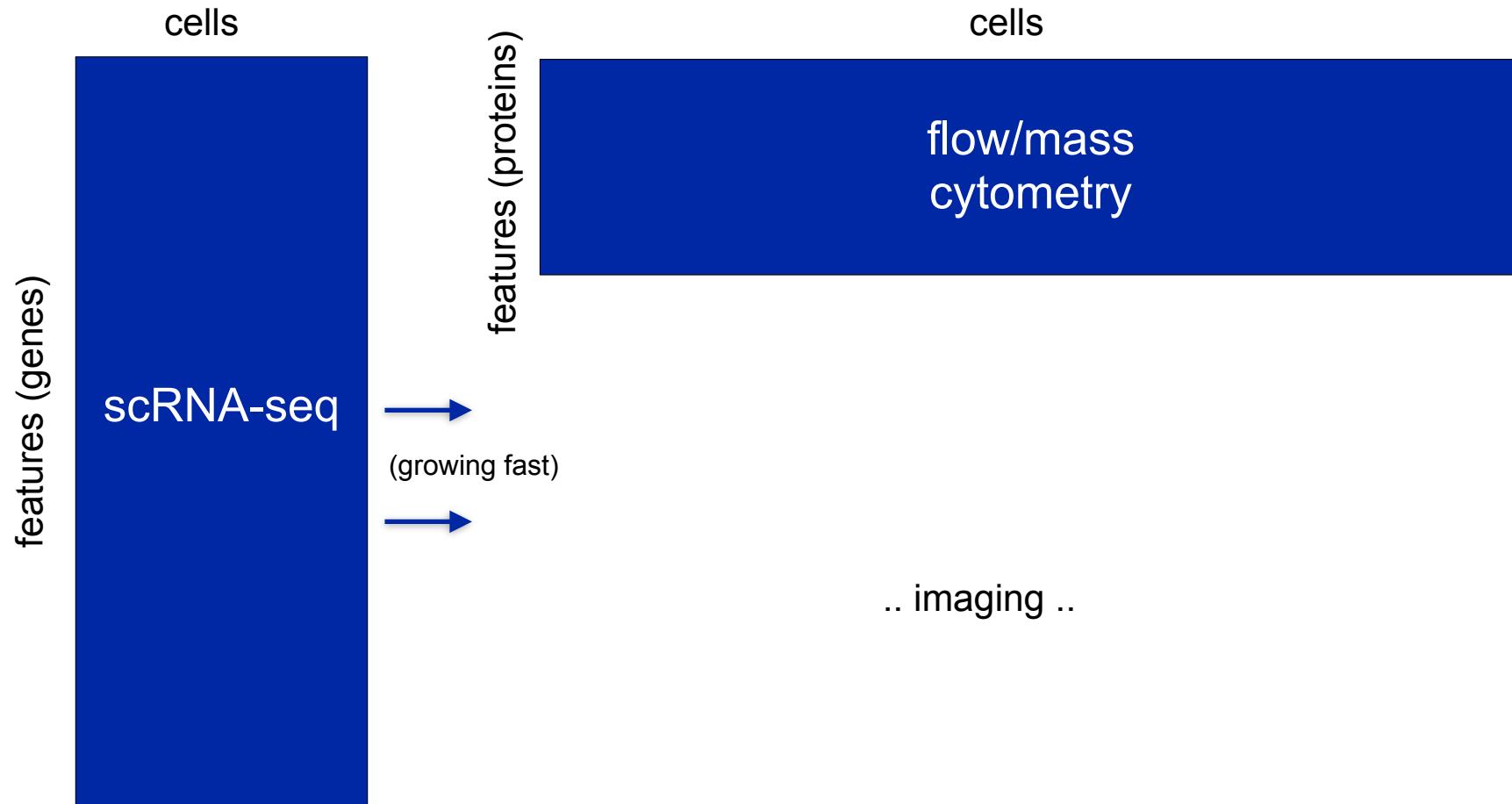
Basic properties: Variability levels



Single-cell RNA-seq: Hypothetical situations



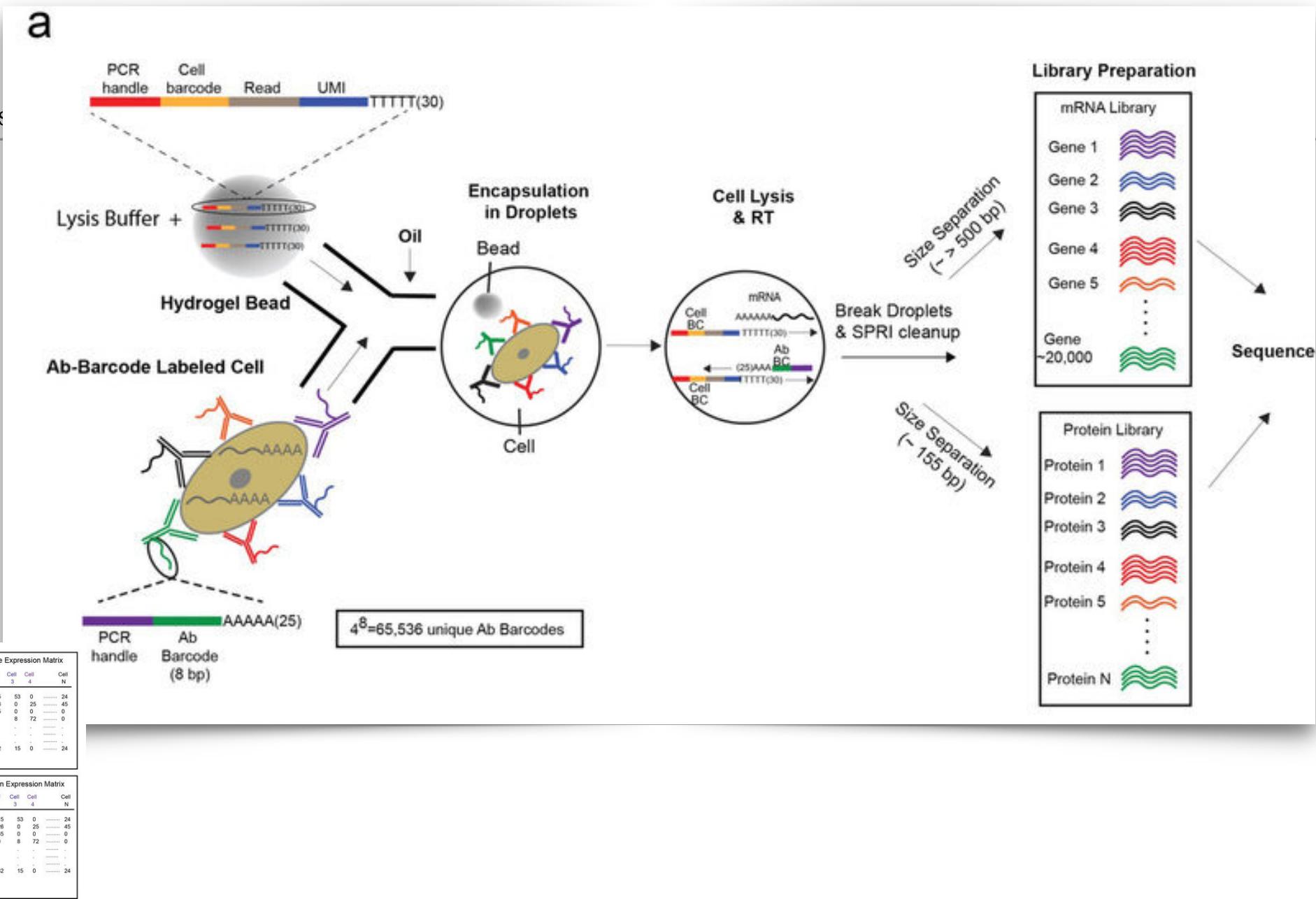
Different shapes of single cell data





Dual assays

- REAP-seq
- CITE-seq





Some terminology: Cell identity, type, state, ..

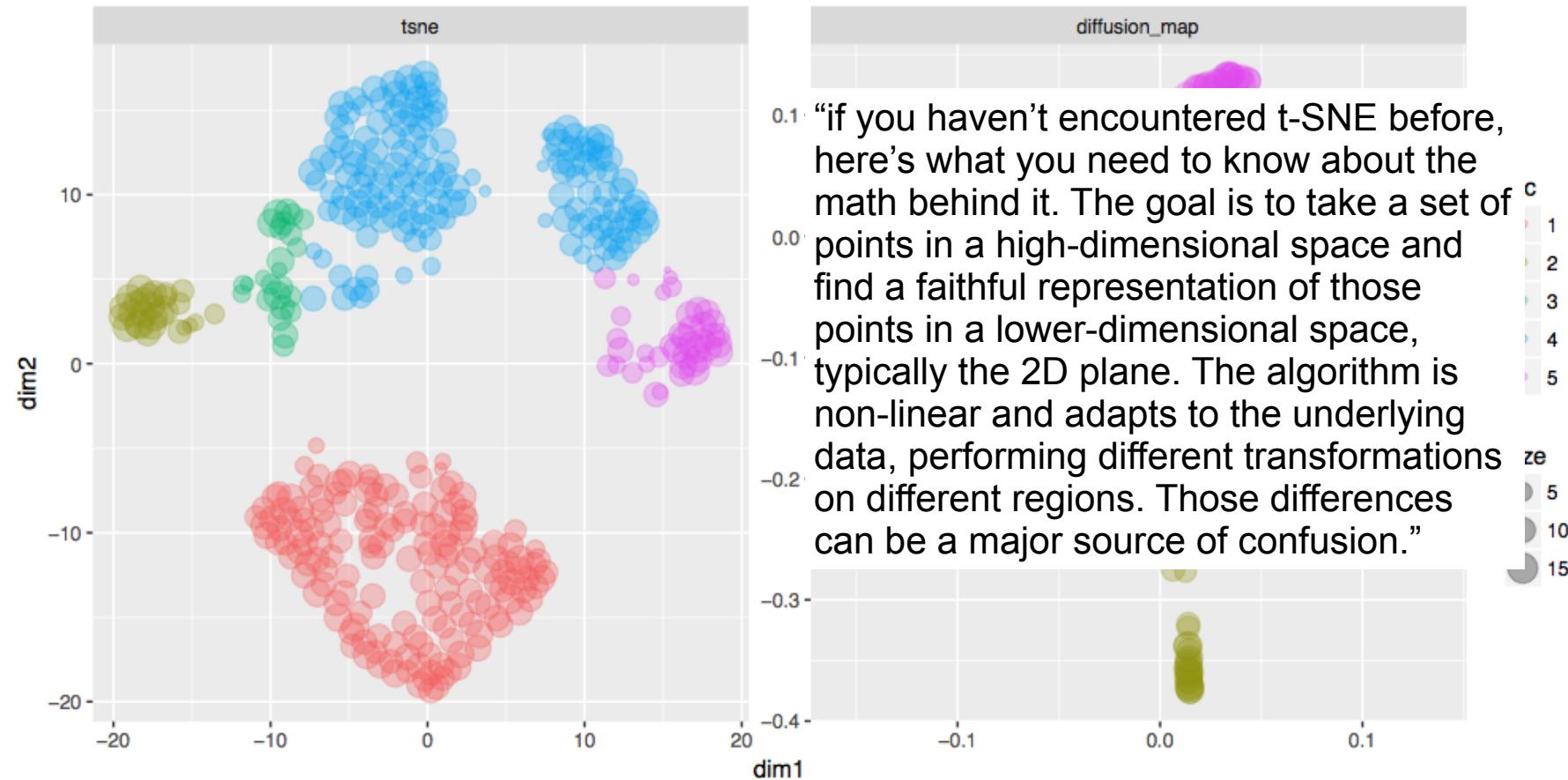
Box 1 The many facets of a cell's identity

We define a cell's identity as the outcome of the instantaneous intersection of all factors that affect it. We refer to the more permanent aspects in a cell's identity as its type (e.g., a hepatocyte typically cannot turn into a neuron) and to the more transient elements as its state. Cell types are often organized in a hierarchical taxonomy, as types may be further divided into finer subtypes; such taxonomies are often related to a cell fate map, reflecting key steps in differentiation. Cell *states* arise transiently during time-dependent processes, either in a *temporal progression* that is unidirectional (e.g., during differentiation, or following an environmental stimulus) or in a *state vacillation* that is not necessarily unidirectional and in which the cell may return to the origin state. Vacillating processes can be *oscillatory* (e.g., cell-cycle or circadian rhythm) or can transition between states with no predefined order (e.g., due to stochastic, or environmentally controlled, molecular events). These time-dependent processes may occur transiently within a stable cell type (as in a transient environmental response), or may lead to a new,

Type: permanent
State: transient



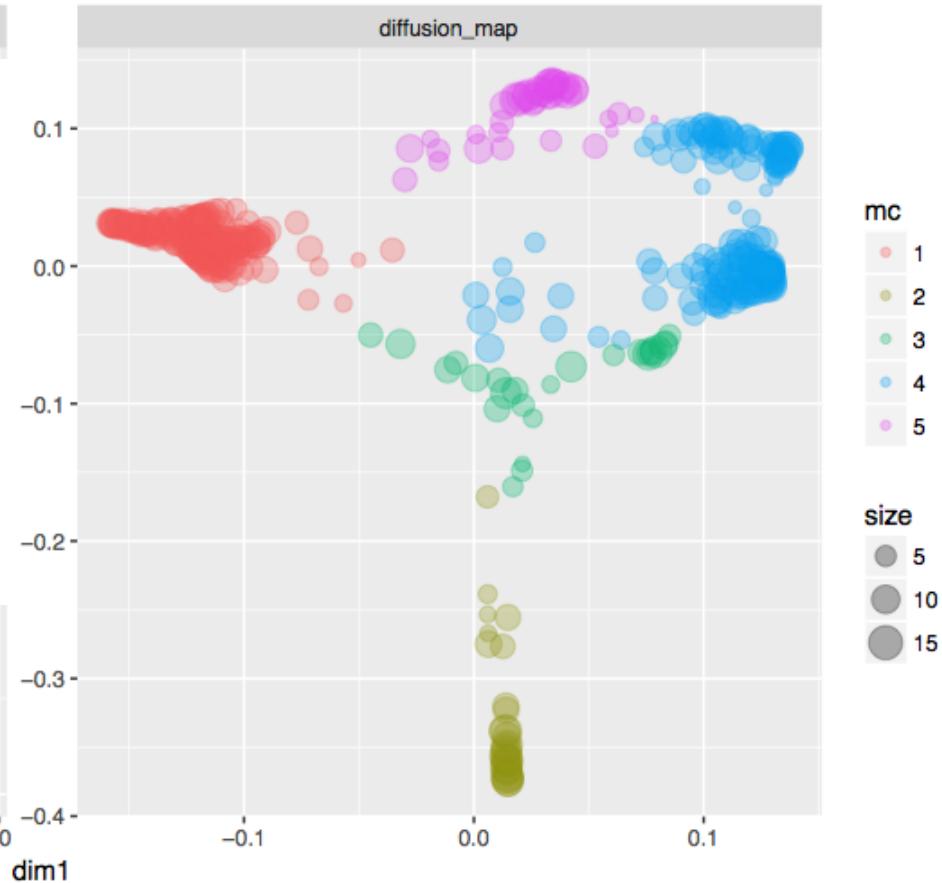
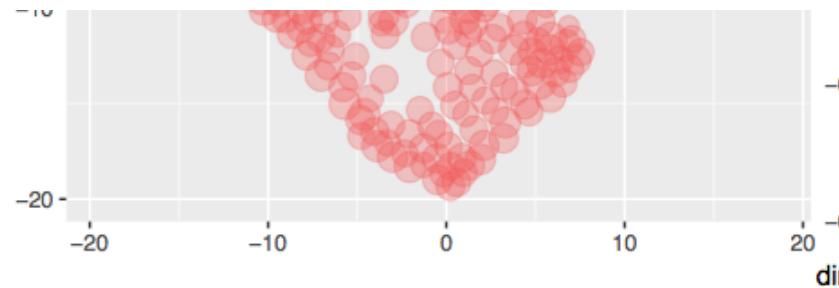
tSNE (t-dist'd stochastic neighbour embedding) + diffusion maps



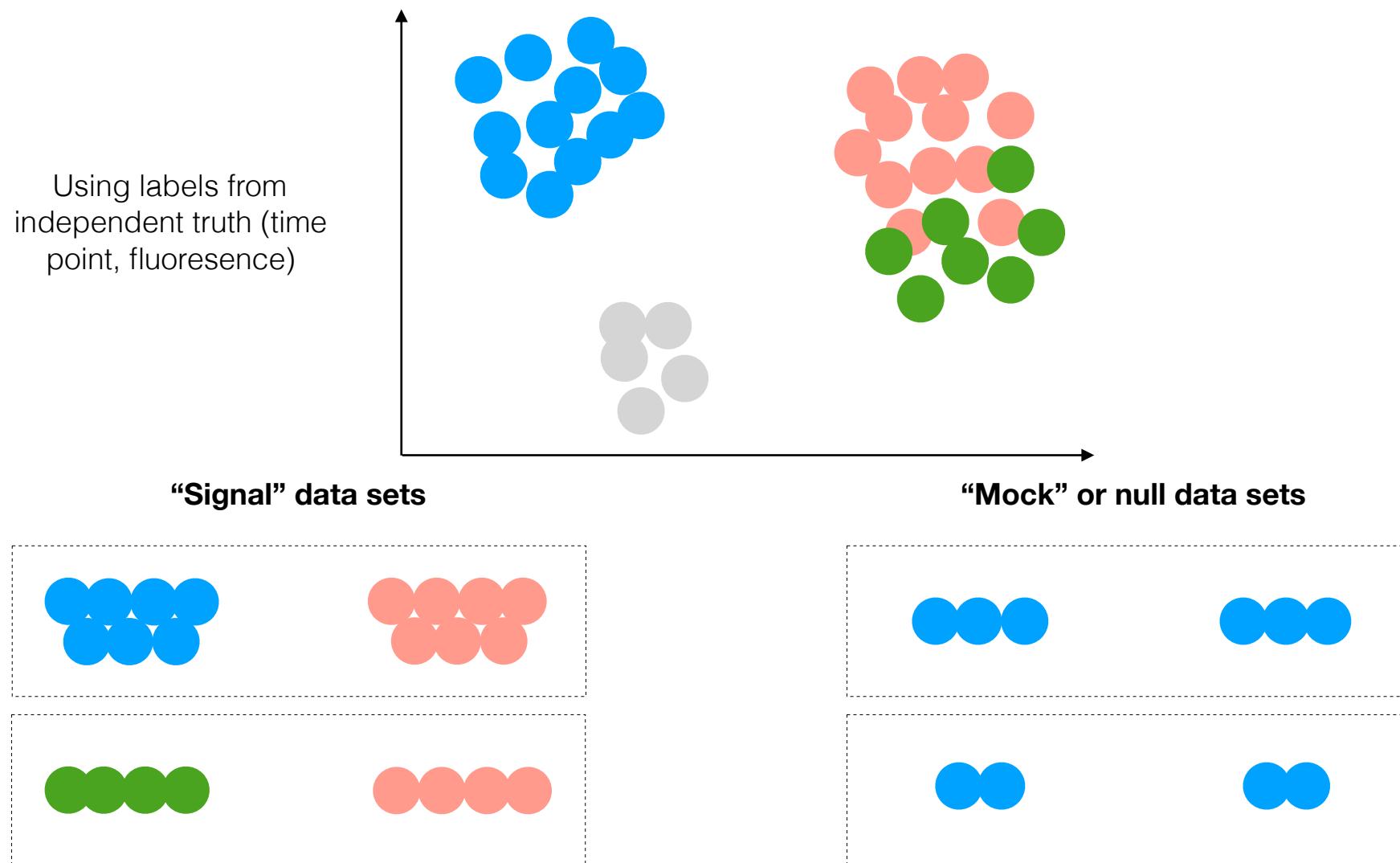


tSNE (t-dist'd stochastic neighbour embedding) + diffusion maps

"Given data in a high-dimensional space .. find parameters that describe the lower-dimensional structures of which it is comprised. Unlike other popular methods such as PCA and MDS, diffusion maps are non-linear and focus on discovering the underlying manifold (lower-dimensional constrained "surface" upon which the data is embedded). By integrating local similarities at different scales, a global description of the data-set is obtained.



Experimental data



Between cell-type DE Benchmark (finding marker genes)

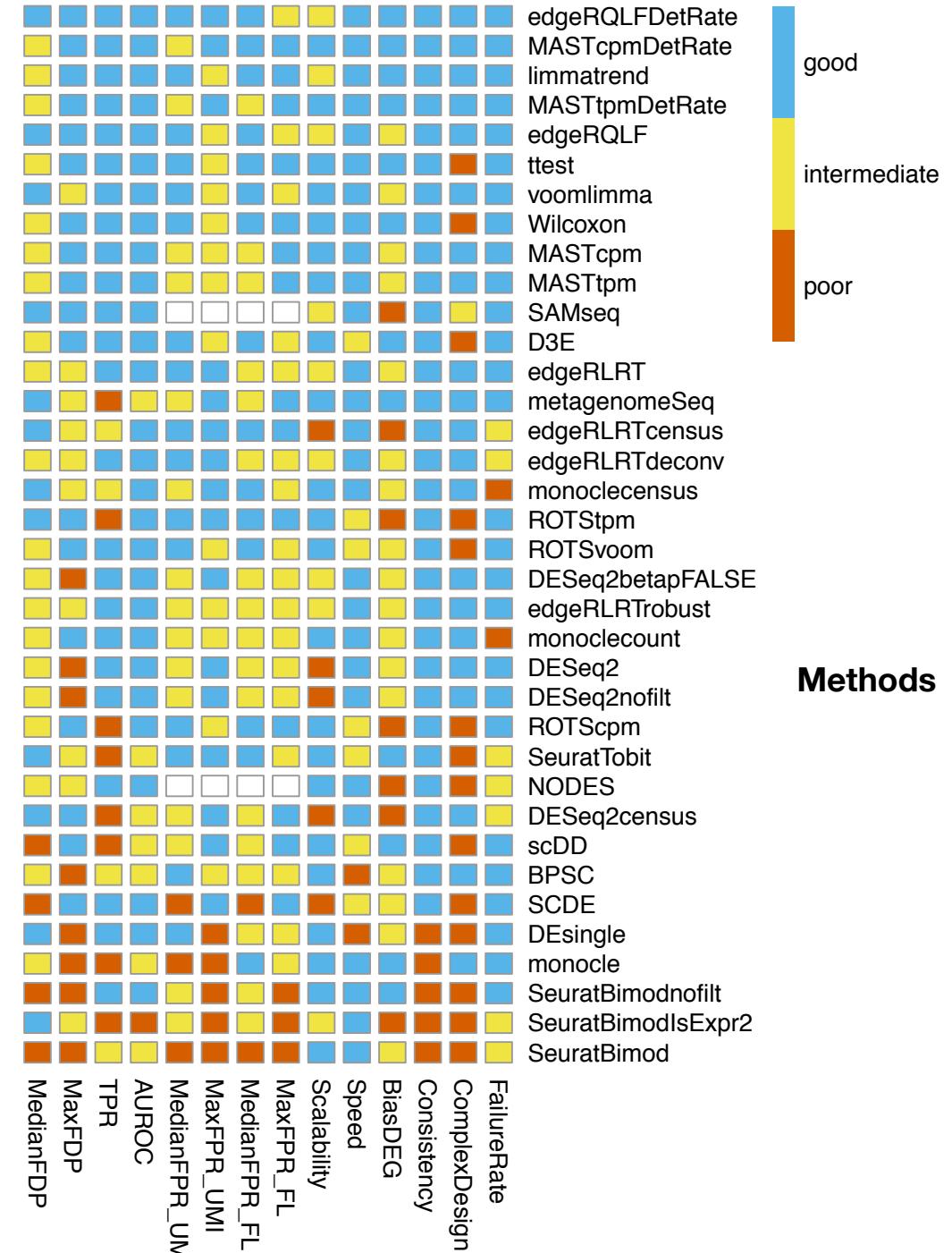
Bias, robustness and scalability in single-cell differential expression analysis

Charlotte Soneson^{1,2}  & Mark D Robinson^{1,2} 

RECEIVED 6 JUNE 2017; ACCEPTED 16 JANUARY 2018; PUBLISHED ONLINE 26 FEBRUARY 2018;

“we found that bulk RNA-seq analysis methods do not generally perform worse than those developed specifically for scRNA-seq.”

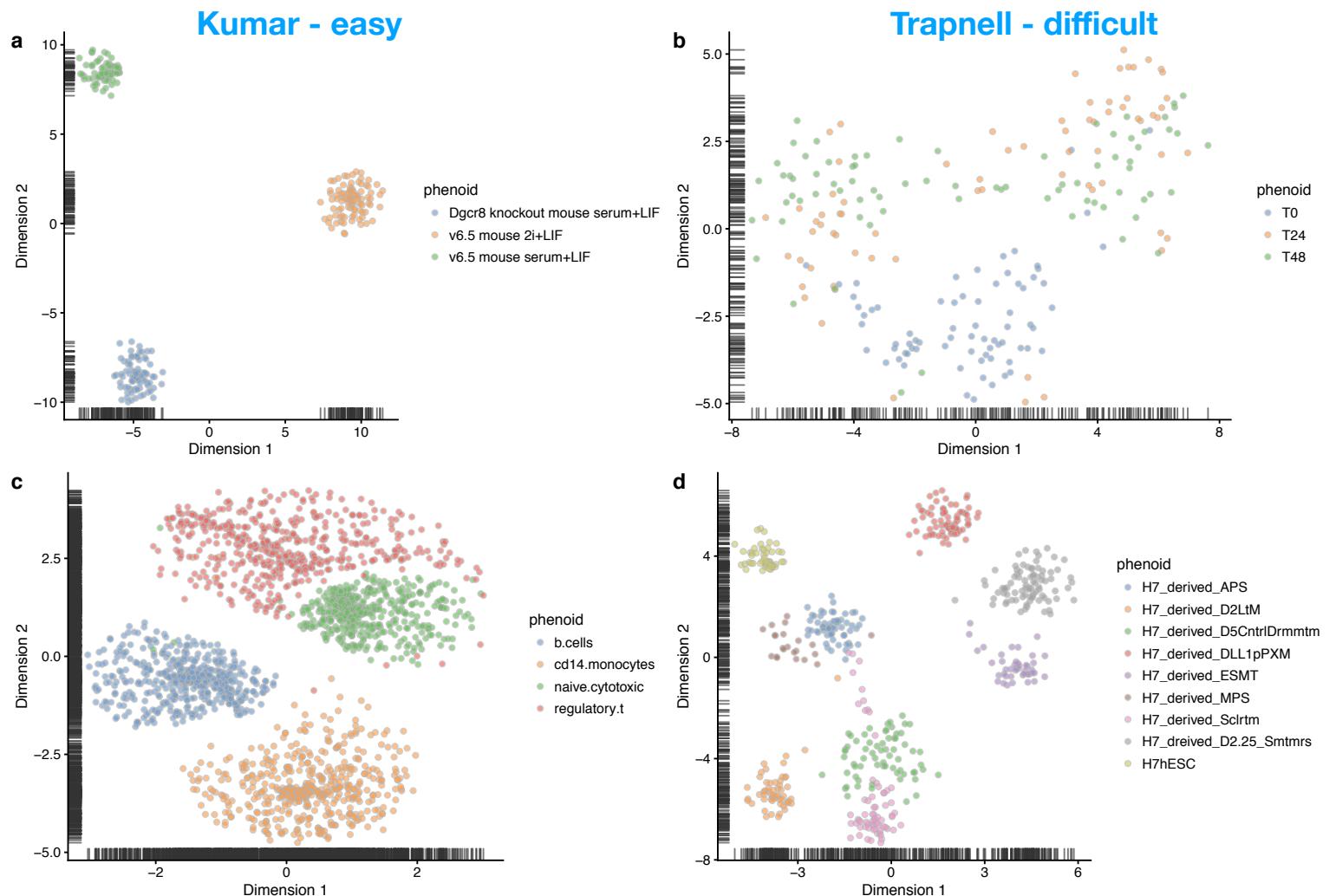
Criteria



Methods

How to cluster scRNA-seq data?

- Datasets from conquer with predefined groups: range of difficulty
- large space: dimension reduction + transformation (log / variance stabilizing) + imputation / zero inflation + clustering method



How to cluster scRNA-seq data?

RESEARCH ARTICLE

REVISED A systematic performance evaluation of clustering methods for single-cell RNA-seq data [version 2; referees: 2 approved]

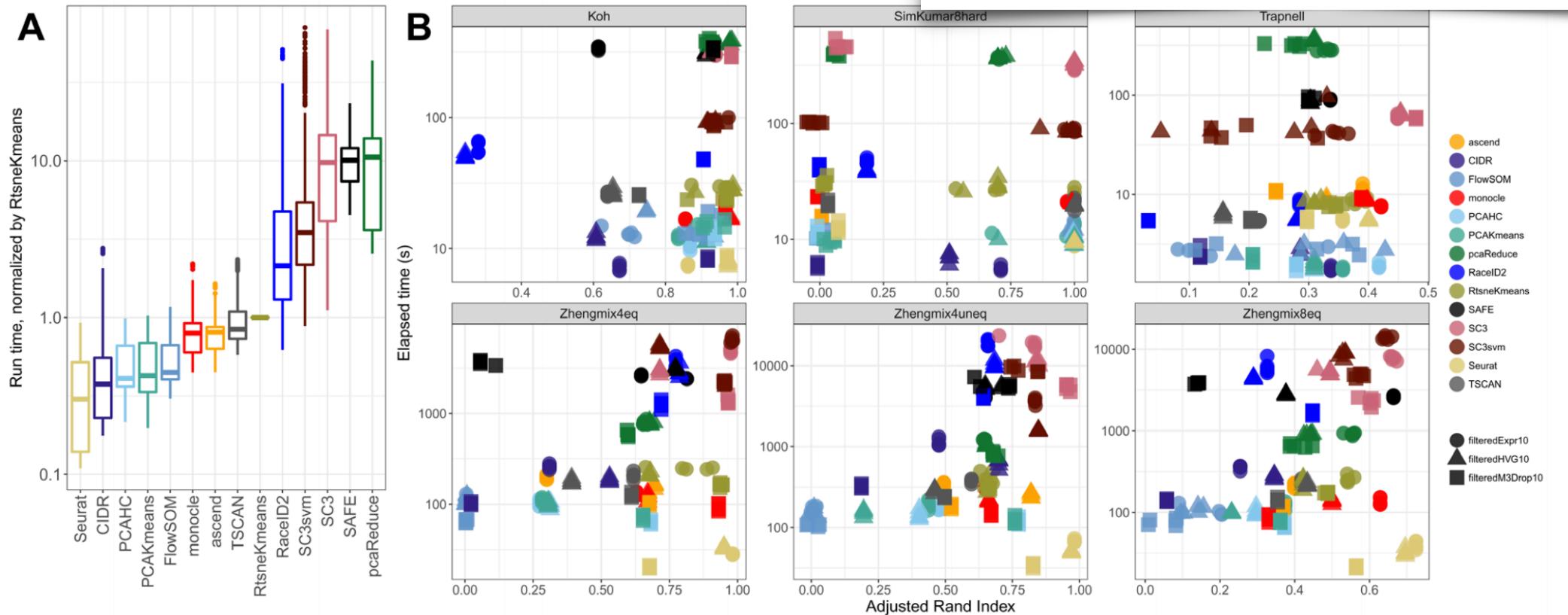
Angelo Duò  1,2, Mark D. Robinson  1,2, Charlotte Soneson  1,2¹Institute of Molecular Life Sciences, University of Zurich, Zurich, 8057, Switzerland²SIB Swiss Institute of Bioinformatics, Zurich, 8057, Switzerland

Figure 2. (A) Normalized run times, using RtsneKmeans as the reference method, across all data set instances and number of clusters. (B) Run time versus performance (ARI) for a subset of data sets and filterings, for the true number of clusters.

conquer - reprocessed data analysis read public scRNA-seq datasets

- <http://imlspenticton.uzh.ch:3838/conquer/>
- Contains both full-length and UMI-based protocols

The screenshot shows a web application titled "conquer" with a blue header bar. Below the header is a navigation menu with tabs: "About", "scRNA-seq data sets" (which is currently selected), "Changelog", "Excluded samples", and "Tutorial". Below the menu, there are search and filter controls: "Show 10 entries" and a "Search:" input field.

The main content area displays a table of scRNA-seq datasets. The columns are labeled: "Data set", "ID", "organism", "ncells", "MultiAssayExperiment", "MultiQC report", "scater report", and "salmon archive". Each dataset row contains a link to the dataset details, the study ID, author, organism, number of cells, and download links for each category.

	Data set	ID	organism	ncells	MultiAssayExperiment	MultiQC report	scater report	salmon archive
1	EMTAB2805 (PMID 25599176)	Buettner2015	Mus musculus	288	Download .rds (2017-07-23)	Download .html (2017-07-23)	Download .html (2017-07-23)	Download .tar.gz (2016-07-23)
2	EMTAB3929 (PMID 27062923)	Petropoulos2016	Homo sapiens	1529	Download .rds (2017-07-22)	Download .html (2017-07-22)	Download .html (2017-07-22)	Download .tar.gz (2017-07-20)
3	GSE41265 (PMID 23685454)	Shalek2013	Mus musculus	18	Download .rds (2017-07-23)	Download .html (2017-07-23)	Download .html (2017-07-23)	Download .tar.gz (2016-09-09)
4	GSE44183-GPL11154 (PMID 23892778)	Xue2013	Homo sapiens	29	Download .rds (2017-07-23)	Download .html (2017-07-23)	Download .html (2017-07-23)	Download .tar.gz (2016-07-24)
5	GSE44183-GPL13112 (PMID 23892778)	Xue2013	Mus musculus	17	Download .rds (2017-07-23)	Download .html (2017-07-23)	Download .html (2017-07-23)	Download .tar.gz (2016-07-24)
6	GSE44183-GPL13112-trimmed (PMID 23892778)	Xue2013	Mus musculus	17	Download .rds (2017-07-23)	Download .html (2017-07-23)	Download .html (2017-07-23)	Download .tar.gz (2016-07-28)



Antonio Rausell

@AntonioRausell

Following



A compilation of recent benchmarks covering each of the prototypical steps in single-cell RNA-seq data analysis. Other references to add here?
cc: @scell_papers #singlecell

Recent benchmarks covering each of the prototypical steps in single-cell RNA-seq data analysis:

Assessment of Single Cell RNA-Seq **Normalization** Methods.
<http://www.g3journal.org/content/7/7/2039.long>

Evaluation of tools for **highly variable gene discovery** from single-cell RNA-seq data
<https://academic.oup.com/bib/advance-article/doi/10.1093/bib/bby011/4898116>

A systematic performance evaluation of **clustering** methods for single-cell RNA-seq data
<https://f1000research.com/articles/7-1141/v2>

Comparison of **clustering** tools in R for medium-sized 10x Genomics single-cell RNA-sequencing data
<https://f1000research.com/articles/7-1297/v1>

Bias, robustness and scalability in single-cell **differential expression** analysis
<https://www.nature.com/articles/nmeth.4612>

A comparison of single-cell **trajectory inference** methods: towards more accurate and robust tools
<https://www.biorxiv.org/content/10.1101/276907v1>

A test metric for assessing single-cell RNA-seq **batch correction**
<https://www.nature.com/articles/s41592-018-0254-1>

scRNA-seq mixology: towards better benchmarking of single cell RNA-seq **protocols** & analysis methods
<https://www.biorxiv.org/content/10.1101/433102v2>

Twitter: @AntonioRausell
antonio.rausell@institutimagine.org