

# **Single-cell sequencing**

# Background

- Most of the biological experiments are performed on “bulk” samples, which contains a large number of cells (millions).
- The high-throughput data we introduced so far are all “bulk” data, which measures the average (gene expression, TF binding, methylation, etc.) of many cells.
- The bulk measurement ignores the inter-cellular heterogeneities:
  - Different cell types.
  - Biological variation among the same type of cell.

# Single-cell biology

- The study of individual cells.
- The cells are isolated from multi-cellular organism.
- Experiment is performed for each cell individually.
- Provides more detailed, higher resolution information.
- High-throughput experiments on single cell is possible.

# Single cell sequencing

- Perform different types of sequencing at the single-cell level:
  - DNA-seq
  - ATAC-seq
  - BS-seq
  - RNA-seq
- Very active research field in the past several years.
- Major challenges:
  - Cell isolation.
  - Amplification of genomic material.
  - Data analysis.

# Basic experimental procedure

- Isolation of single cell. Techniques include
  - Laser-capture microdissection (LCM)
  - Fluorescence-activated cell sorting (FACS)
  - Microfluidics
- Open the cell and obtain DNA/mRNA/etc.
- PCR amplification to get enough materials.
- Perform sequencing.
- Note that single cell sequencing usually has higher error rates than bulk data.

# Single cell DNA-seq (scDNA-seq)

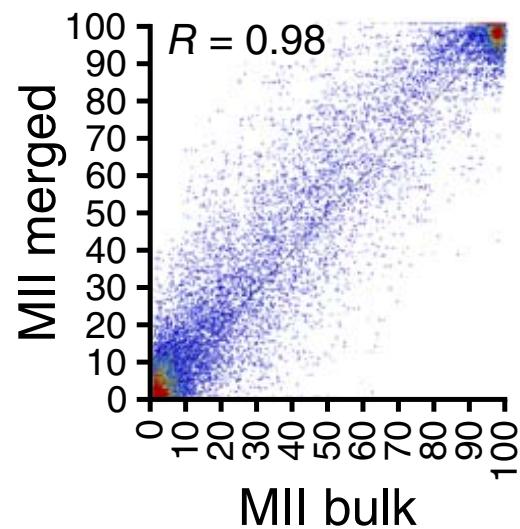
- For a comprehensive review, read *Gawad et al.* (2016) NRG.
- Examples of biological applications:
  - Identify and assemble the genome of unculturable microorganisms.
  - Determine the contribution of intra-tumor genetic heterogeneity in cancer development of treatment response.

# scDNA-seq data analysis

- Single cell variant calling:
  - Bulk data can be used as reference to reduce false positives.
  - Combine data from several cells.
  - Software: Monovar (*Zafar et al. 2016 Nat. Method.*)
- Determining genetic relationship among single cells:
  - This is a clustering problem. Cells can be put into groups or a phylogenetic tree based on similarity of variants.
  - Methods are mostly ad hoc.

# Single cell BS-seq (scBS-seq)

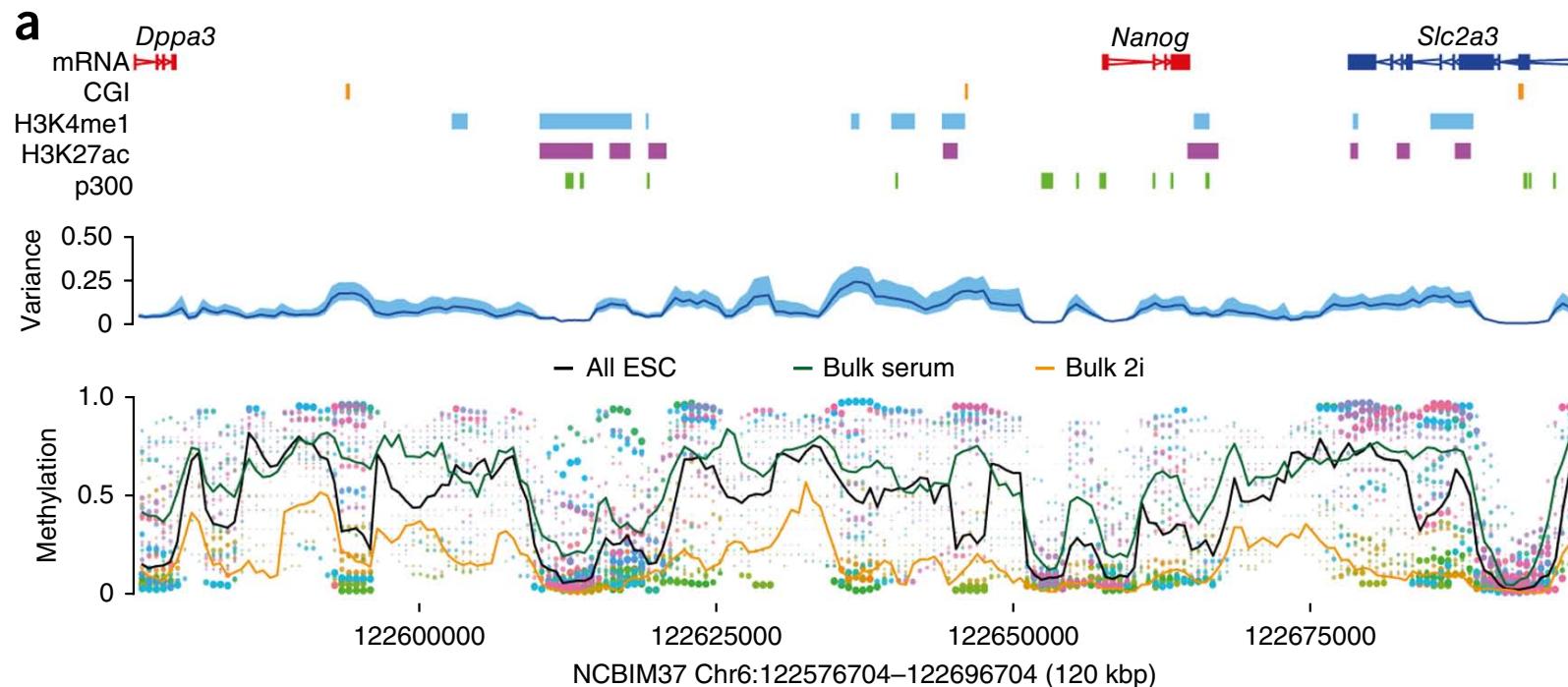
- Similar to scDNA-seq, but with bisulfite treatment before sequencing.
- There's scWGBS and scRRBS.
- The methylation levels from scBS-seq should be 0/1, with some exceptions caused by technical artifacts.
- Merged single cell and bulk data have good correlation.

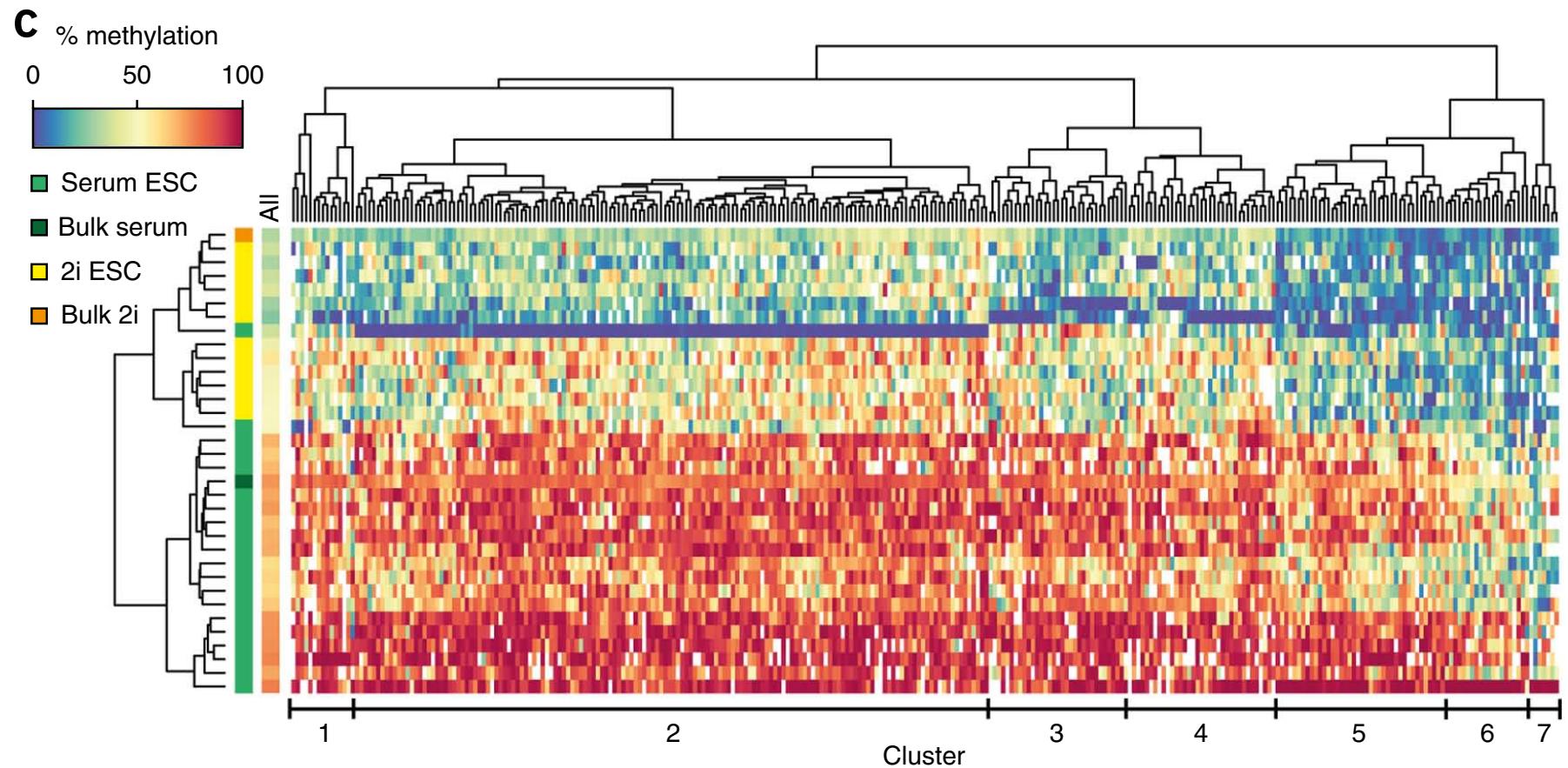


Smallwood et al. 2014, NM

# scBS-seq data analysis

- So far the data analysis are mostly descriptive:
  - compute variations among cells
  - Cell clustering



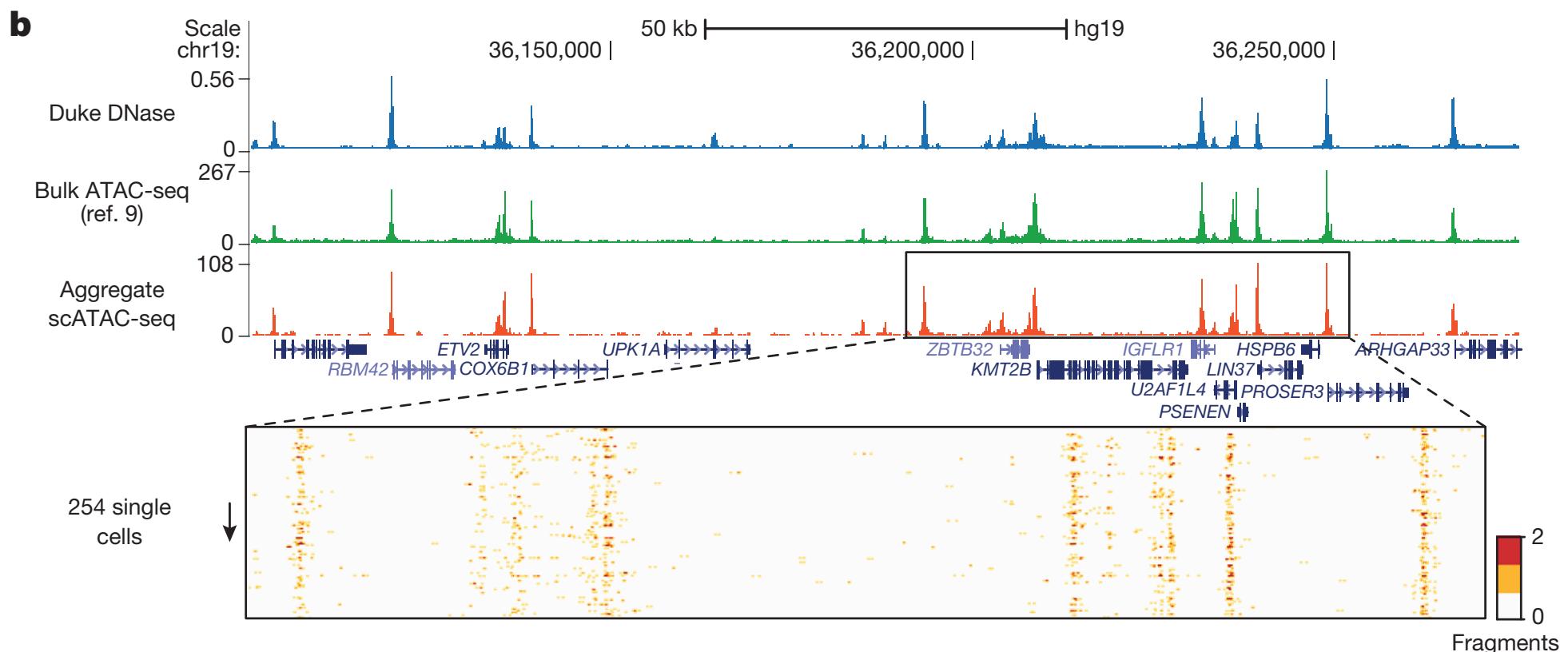


# Single cell ChIP/ATAC-seq

- ATAC-seq: similar to DNase-seq, profile the active genomic regions. Data look like ChIP-seq.
- A few papers:
  - Rotem et al. (2015) NBT: scChIP-seq
  - Buenrostro et al. (2015) Nature: scATAC-seq
  - Pott and Liet (2015) Genome Biology: review

# scChIP/scATAC-seq data

- Aggregated sc data has good agreement with bulk.



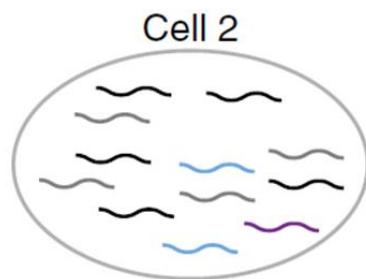
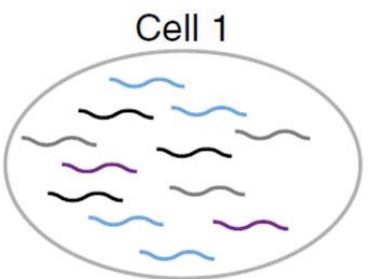
- Very sparse: one or a few reads at peak regions.
  - Extremely low signal to noise ratio.
  - Peak calling have to be based on combined data, or rely on other prior information

	Peak1	Peak2	Peak3	...													
Cell1	1	1	2	1	0	2	2	1	0	2	1	0	1	0	0	0	2
Cell2	0	1	0	1	0	0	1	0	0	0	0	3	1	0	1	1	0
Cell3	2	0	0	1	1	0	1	0	1	0	1	1	1	2	1	0	0
⋮	1	1	2	1	0	2	2	1	0	2	1	0	1	0	0	0	2
⋮	0	1	0	1	0	0	1	0	0	0	0	3	1	0	1	1	0
	2	0	0	1	1	0	1	0	1	0	1	1	1	2	1	0	0
	1	1	2	1	0	2	2	1	0	2	1	0	1	0	0	0	2

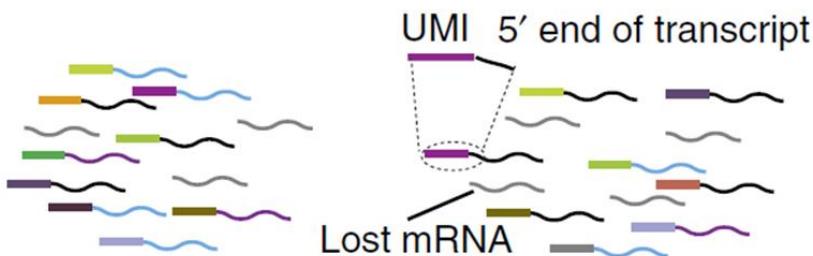
# Single cell RNA-seq (scRNA-seq)

- The most active in the sc field.
- Scientific goals:
  - Understand the gene expression heterogeneity within the same sample.
  - Composition of different types of cell in complex tissues, such as brain, cancer, etc.
  - Above can be explored spatially, temporally, or under different biological condition.
- Raw data are the same as bulk RNA-seq, can be aligned using the same software.

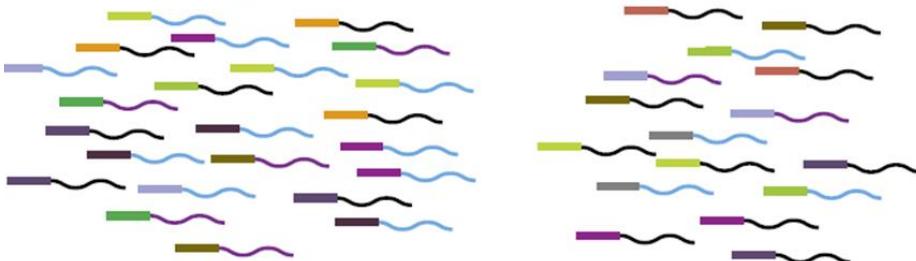
# Experimental procedure



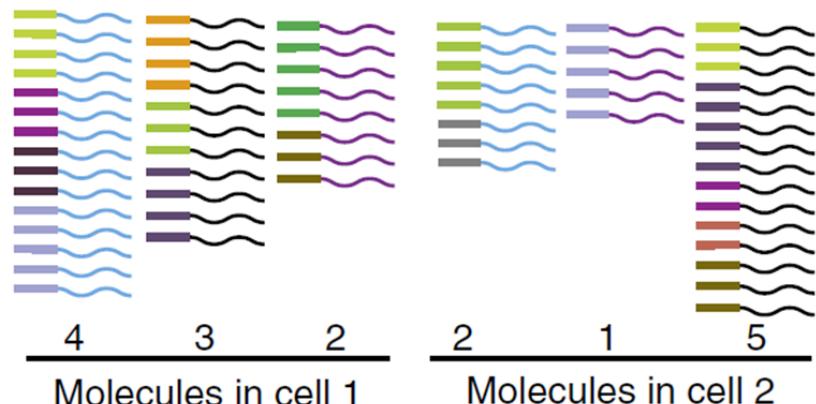
Reverse transcription, barcoding and UMI labeling



PCR amplification



Sequencing and computation



Saiful Islam ... Sten Linnarsson

Cell barcode      UMI      cDNA (50-bp sequenced)

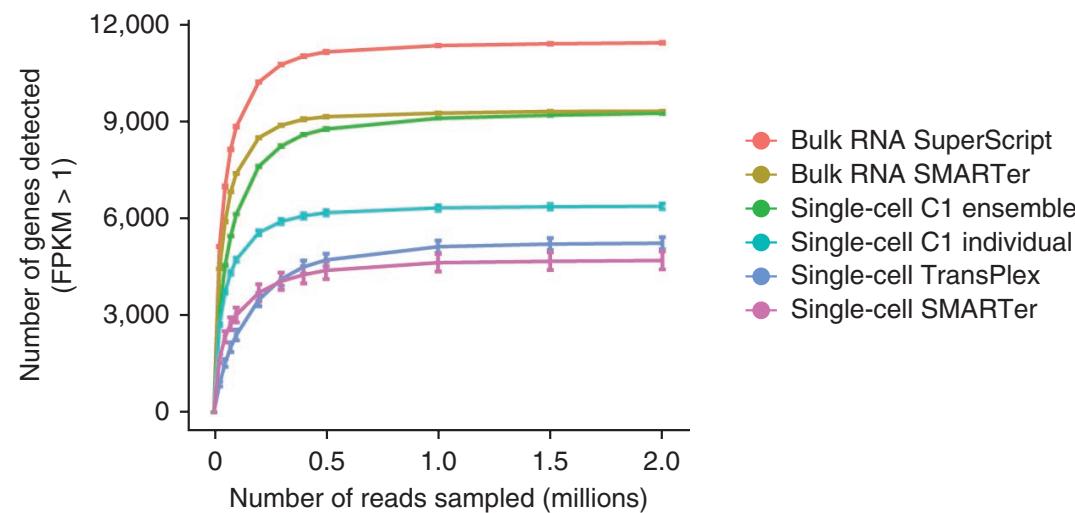
AAATTATGACGA	TGTGCTTG	..... GACTGCAC
CGTTAGATGGCA	GGGCCGGG	..... CTCATAGT
GACCTACGAGTT	AGTTTGTA	..... GCTCATAA
GTAAACGTACC	CTAGCTGT	..... GATTTTCT
ACGTCACCTTT	GTGGGGGT	..... ATAAGCTC
TTGCCGTGGTGT	TATGGAGG	..... CCAGCACC
AGTCCATGTGCGGCAGGTT	..... GTTGGCGT	
AAATTATGACGA	AGTTTGTA	..... AGATGGGG
CCAAAGATGTCC	TCTAGGCT	..... GGGGACGA
GTAAACGTACC	AAGGCTTG	..... CAAAGTTC
TTTTGACCAGT	CGTGAGGG	..... TTCCAAGG
ACTGTCCATGCC	CCTGTGTA	..... TGTTACGT
CGTAAAACAATA	ATCCGGTG	..... TTAAACCG
.....	.....	.....

cDNA alignment to genome and group results by cell

Cell 1	{	TTGCCGTGGTGT	GGCGGGGA..... CGGTGTTA ]	<i>DDX51</i>
		TTGCCGTGGTGT	TATGGAGG..... CCAGCACC ]	<i>NOP2</i>
		TTGCCGTGGTGT	TCTCAAGT..... AAAATGGC ]	<i>ACTB</i>
Cell 2	{	CGTTAGATGGCA	GGGCCGGG..... CTCATAGT ]	<i>LBR</i>
		CGTTAGATGGCA	ACGTTATA..... ACGCGTAC ]	<i>ODF2</i>
		CGTTAGATGGCA	TCGAGATT..... AGCCCTTT ]	<i>HIF1A</i>
Cell 3	{	AAATTATGACGA	AGTTTGTA..... GGGAAATTA ]	<i>ACTB</i>
		AAATTATGACGA	AGTTTGTA..... AGATGGGG ]	
		AAATTATGACGA	TGTGCTTG..... GACTGCAC ]	<i>RPS15</i>
Cell 4	{	GTAAACGTACC	CTAGCTGT..... GATTTTCT ]	<i>GTPBP4</i>
		GTAAACGTACC	GCAGAACT..... GTTGGCGT ]	<i>GAPDH</i>
		GTAAACGTACC	AAGGCTTG..... CAAAGTTC ]	
		GTAAACGTACC	TTCCGGTC..... TCCAGTCG ]	<i>ARL1</i>

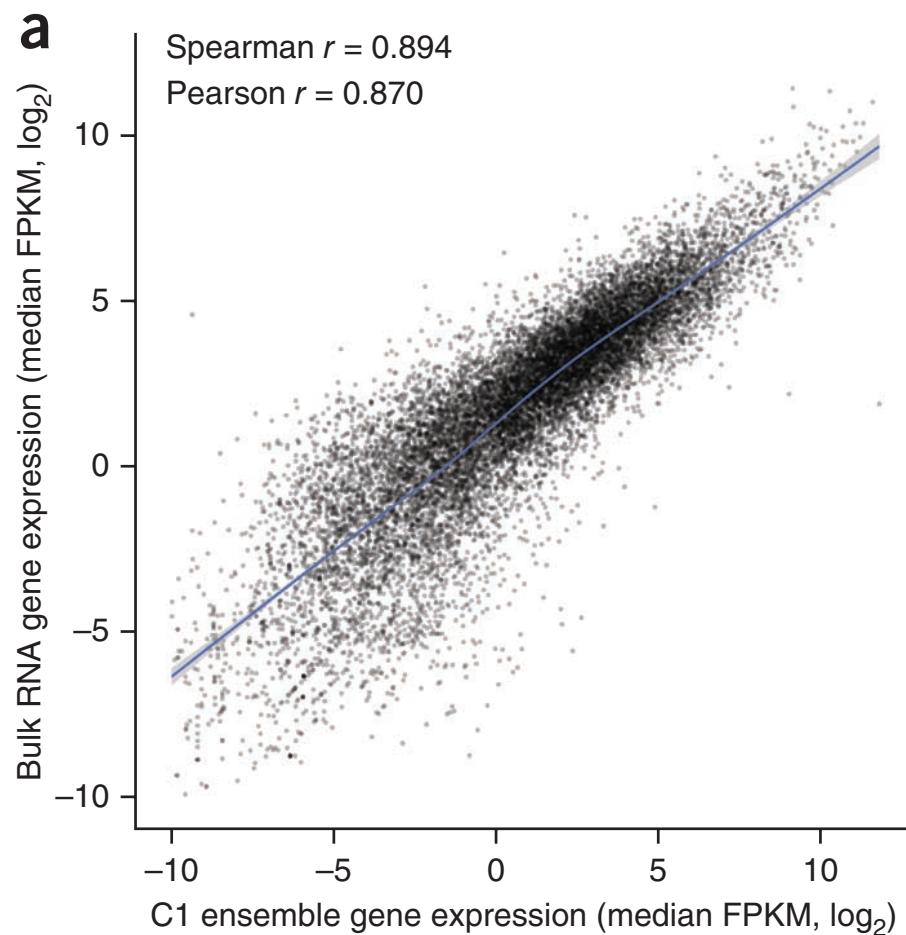
# Some data characteristics

- Number of transcripts detected is much lower compared to bulk RNA-seq, due to low capture and reverse transcription efficiencies.

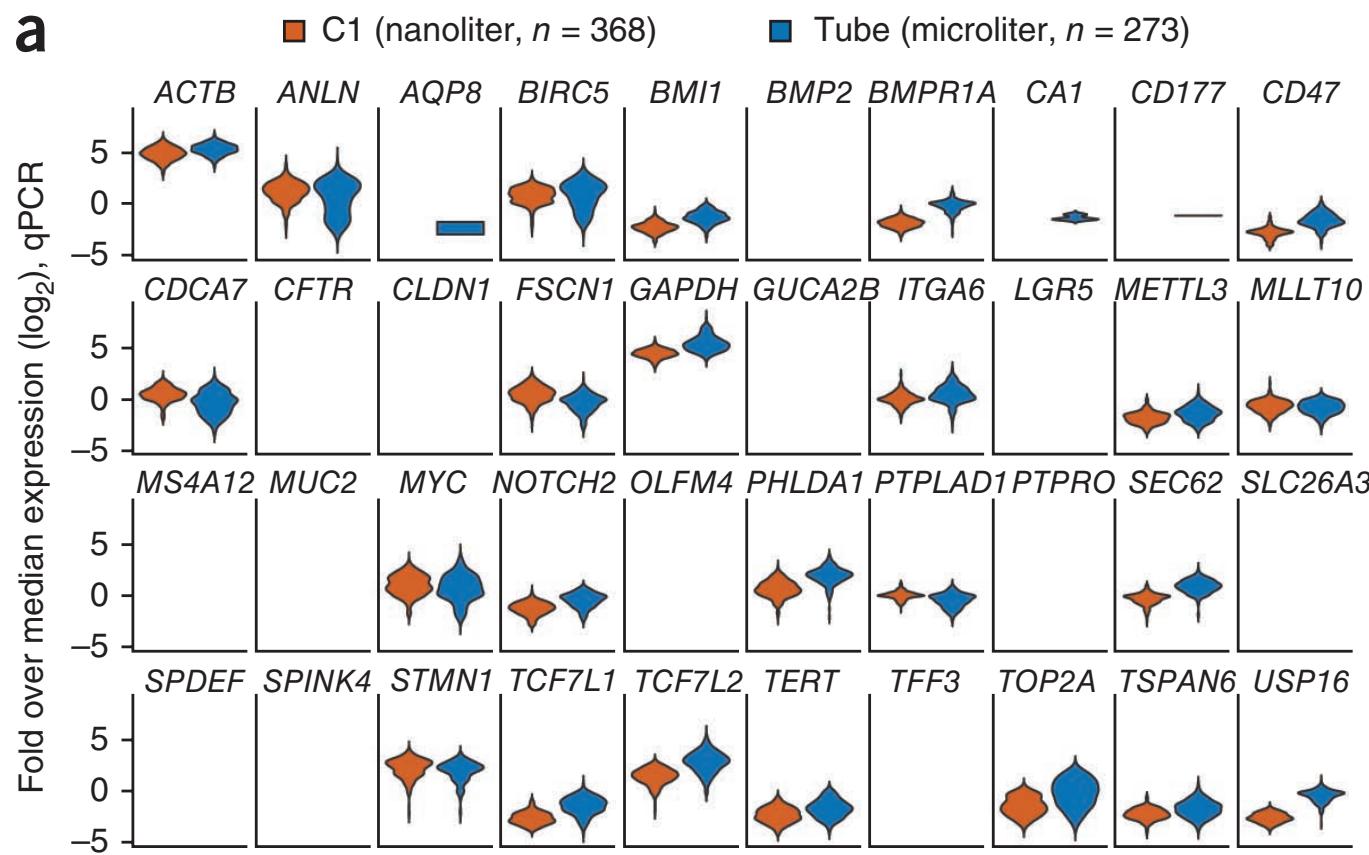


**Figure 5 |** Saturation curves for the different sample preparation methods. Each point on the curve was generated by randomly selecting a number of raw reads from each sample library and then using the same alignment pipeline to call genes with mean FPKM >1. Each point represents four replicate subsamplings. Error bars, standard error.

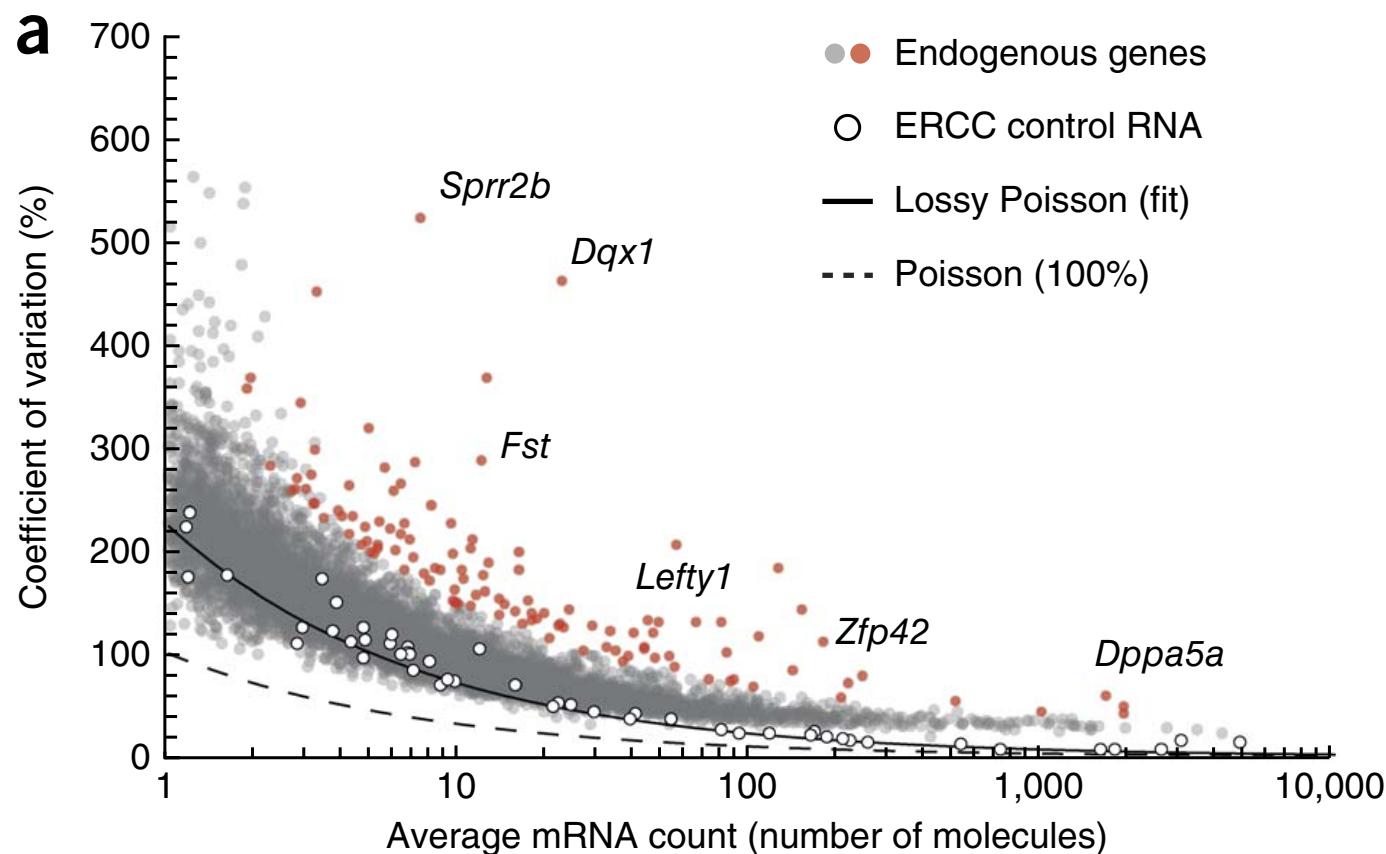
- Bulk and aggregated single cell expressions have good correlation.



- Expression levels for a gene in different cells sometimes show bimodal distribution.



- Negative correlation between mean expression and biological variation (same as in bulk).



# Normalization issues

- scRNA-seq is very noisy.
- Spike-in data is usually available.
  - Spike-ins from the external RNA Control Consortium (ERCC) panel, which contains 92 synthetic spikes based on bacterial genome.
- UMI (unique molecule identifier) is sometimes used to barcode the molecules for estimating amplification noise.
- A combination of spike-in and UMI can potentially be used for data normalization.

# Existing work for scRNA-seq normalization

---

*Application Note*

## Normalization and noise reduction for single cell RNA-seq experiments

Bo Ding<sup>1,#</sup>, Lina Zheng<sup>1,#</sup>, Yun Zhu<sup>1</sup>, Nan Li<sup>1</sup>, Haiyang Jia<sup>1,2</sup>, Rizi Ai<sup>1</sup>, Andre Wildberg<sup>1</sup> and Wei Wang<sup>1,3\*</sup>

<sup>1</sup>Department of Chemistry and Biochemistry, University of California, La Jolla, CA 92093, USA,

<sup>2</sup> College of Computer Science and Technology, Jilin University, Changchun 130012, China.

<sup>3</sup>Department of Cellular and Molecular Medicine, University of California, La Jolla, CA 92093, USA,

#Equal contribution

Associate Editor: Dr. Ziv Bar-Joseph

---

- Log-transform FPKM values, denoted by  $x$ .
- Assume the expression value,  $y$ , follow Gamma distribution. The mean of Gamma is a polynomial function of  $x$ :  $y = \mu(x)$ .

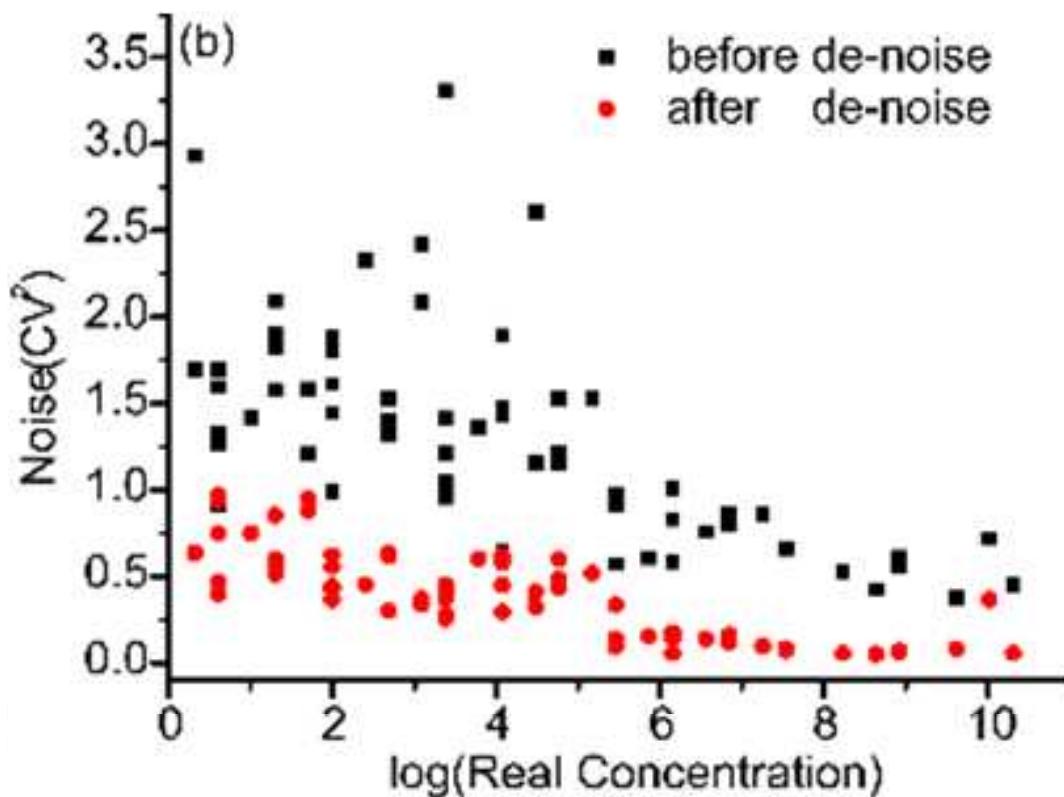
$$\mu(x) = \sum_{i=0}^n \beta_i x^i.$$

The model is the following:

$$y \sim \text{Gamma}(y; \mu(x), \varphi)$$

- Use MLE to estimate parameters based on ERCC data. Then the fitted model is applied to all genes to estimate concentration.

- Results: reduced CV cross cells.



METHOD

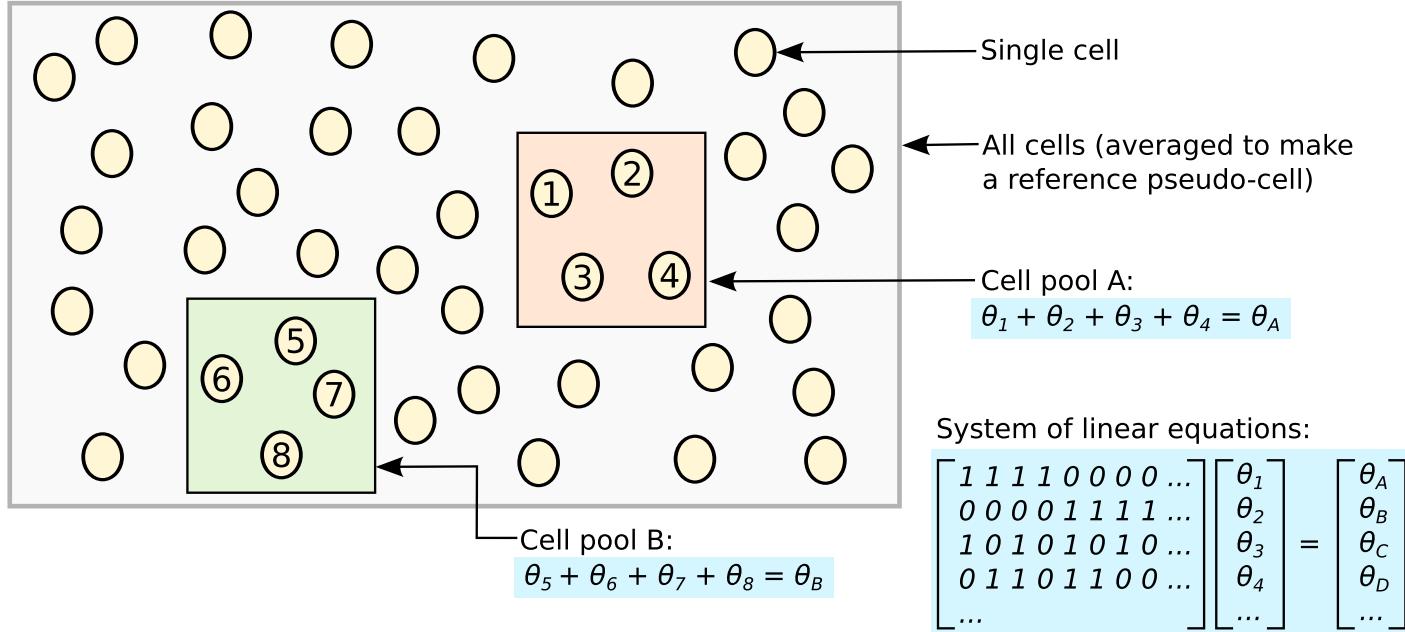
Open Access



# Pooling across cells to normalize single-cell RNA sequencing data with many zero counts

Aaron T. L. Lun<sup>1\*</sup>, Karsten Bach<sup>2</sup> and John C. Marioni<sup>1,2,3\*</sup>

- Works for data without spike-in.
- The goal is to estimate a size factor for each cell.
- The idea is to normalize on summed expression values from pools of cells – it's more stable than using individual cell.



**Fig. 3** Schematic of the deconvolution method. All cells in the data set are averaged to make a reference pseudo-cell. Expression values for cells in pool A are summed together and normalized against the reference to yield a pool-based size factor  $\theta_A$ . This is equal to the sum of the cell-based factors  $\theta_j$  for cells  $j = 1-4$  and can be used to formulate a linear equation. (For simplicity, the  $t_j$  term is assumed to be unity here.) Repeating this for multiple pools (e.g., pool B) leads to the construction of a linear system that can be solved to estimate  $\theta_j$  for each cell  $j$

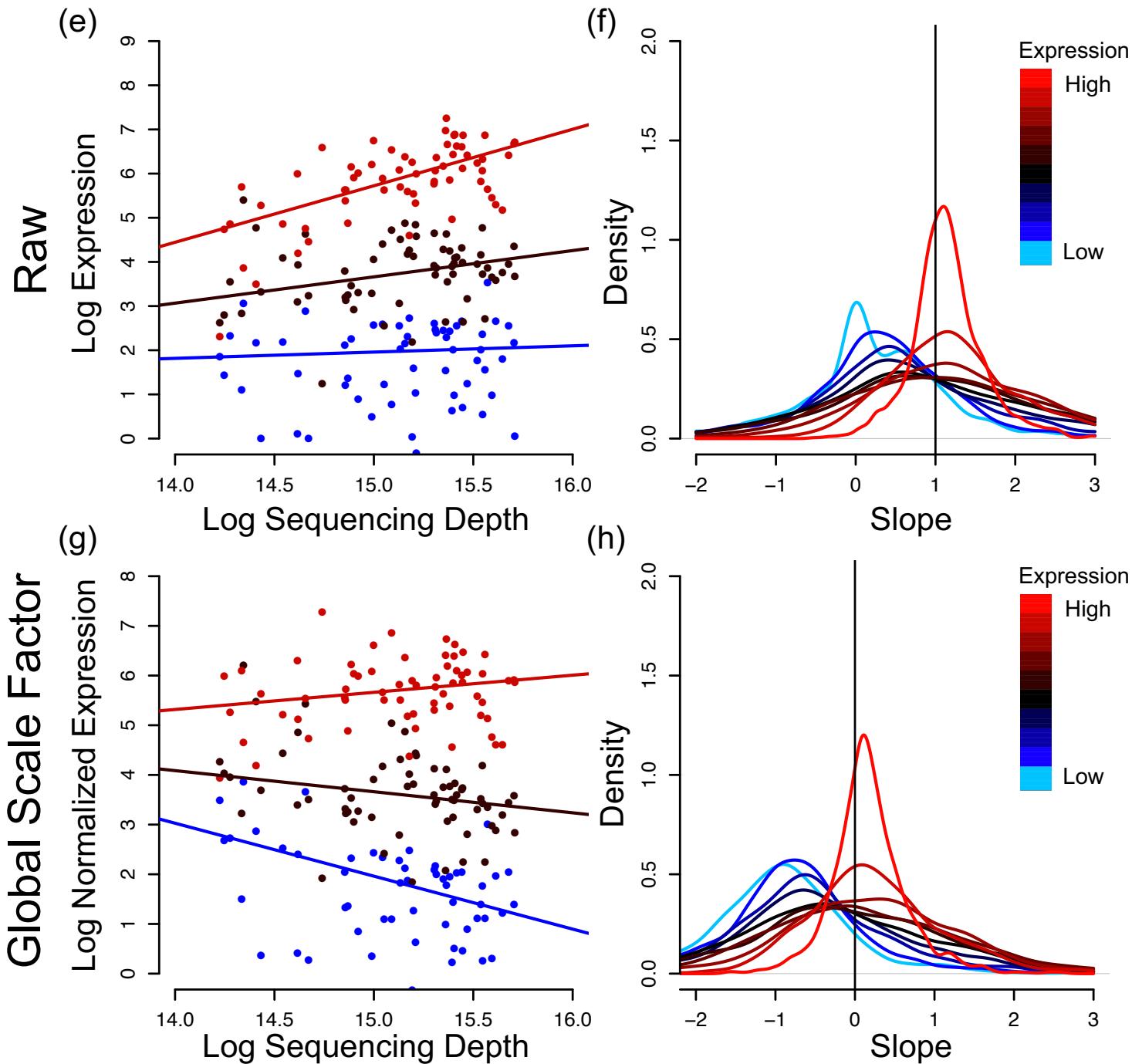
# SCnorm: robust normalization of single-cell RNA-seq data

Rhonda Bacher<sup>1,5</sup> , Li-Fang Chu<sup>2,5</sup>, Ning Leng<sup>2</sup>,  
Audrey P Gasch<sup>3</sup>, James A Thomson<sup>2</sup>, Ron M Stewart<sup>2</sup>,  
Michael Newton<sup>1,4</sup>  & Christina Kendziorski<sup>4</sup>

584 | VOL.14 NO.6 | JUNE 2017 | NATURE METHODS

- Basic idea: one normalization factor per cell doesn't fit all genes.
- Relationships of read counts and sequencing depths vary and depend on the expression levels.

# Single cell



# scNorm procedure

- Uses quantile regression to estimate the dependence of read counts on sequencing depth for every gene.
- Genes with similar dependence are then grouped, and a second quantile regression is used to estimate scale factors within each group.
- Implemented in software SCnorm.

# Differential expression

- Traditional methods test mean changes.
- Due to the bi-modal distribution of the GE in scRNA-seq, the consideration and modeling of “drop-out” event (non-expressed) is very important.
- A long list of tools: MAST, SCDE, BPSC, DEsingle, Seurat, Monocle, SC2P, scDD, ...

# Bayesian approach to single-cell differential expression analysis

740 | VOL.11 NO.7 | JULY 2014 | NATURE METHODS

Peter V Kharchenko<sup>1-3</sup>, Lev Silberstein<sup>3-5</sup> &  
David T Scadden<sup>3-5</sup>

- SCDE (single-cell differential expression).
- Use a mixture of a Poisson with small rate (dropout) and negative binomial (expressed) to model the expression:  $p(x | r_c, \Omega_c) = p_d(x)p_{Poisson}(x) + (1 - p_d(x))p_{NB}(x | r_c)$
- The DE is based on Bayesian inference. But the derivation in this paper is messy.

METHOD

Open Access



# MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data

Greg Finak<sup>1†</sup>, Andrew McDavid<sup>1†</sup>, Masanao Yajima<sup>1†</sup>, Jingyuan Deng<sup>1</sup>, Vivian Gersuk<sup>2</sup>, Alex K. Shalek<sup>3,4,5,6</sup>, Chloe K. Slichter<sup>1</sup>, Hannah W. Miller<sup>1</sup>, M. Juliana McElrath<sup>1</sup>, Martin Prlic<sup>1</sup>, Peter S. Linsley<sup>2</sup>  
and Raphael Gottardo<sup>1,7\*</sup>

- MAST: “Model-based Analysis of Single- cell Transcriptomics.”

# MAST for DE

- Main ideas:
  - Use  $\log_2(\text{TPM}+1)$  as input data
  - Both dropout probability and expression level depends on experimental conditions.

$$\text{logit}(\Pr(Z_{ig} = 1)) = X_i \beta_g^D$$

$$\Pr(Y_{ig} = y | Z_{ig} = 1) = N(X_i \beta_g^C, \sigma_g^2)$$

- Model fitting with some regularization.
- DE is based on chi-square or Wald test.

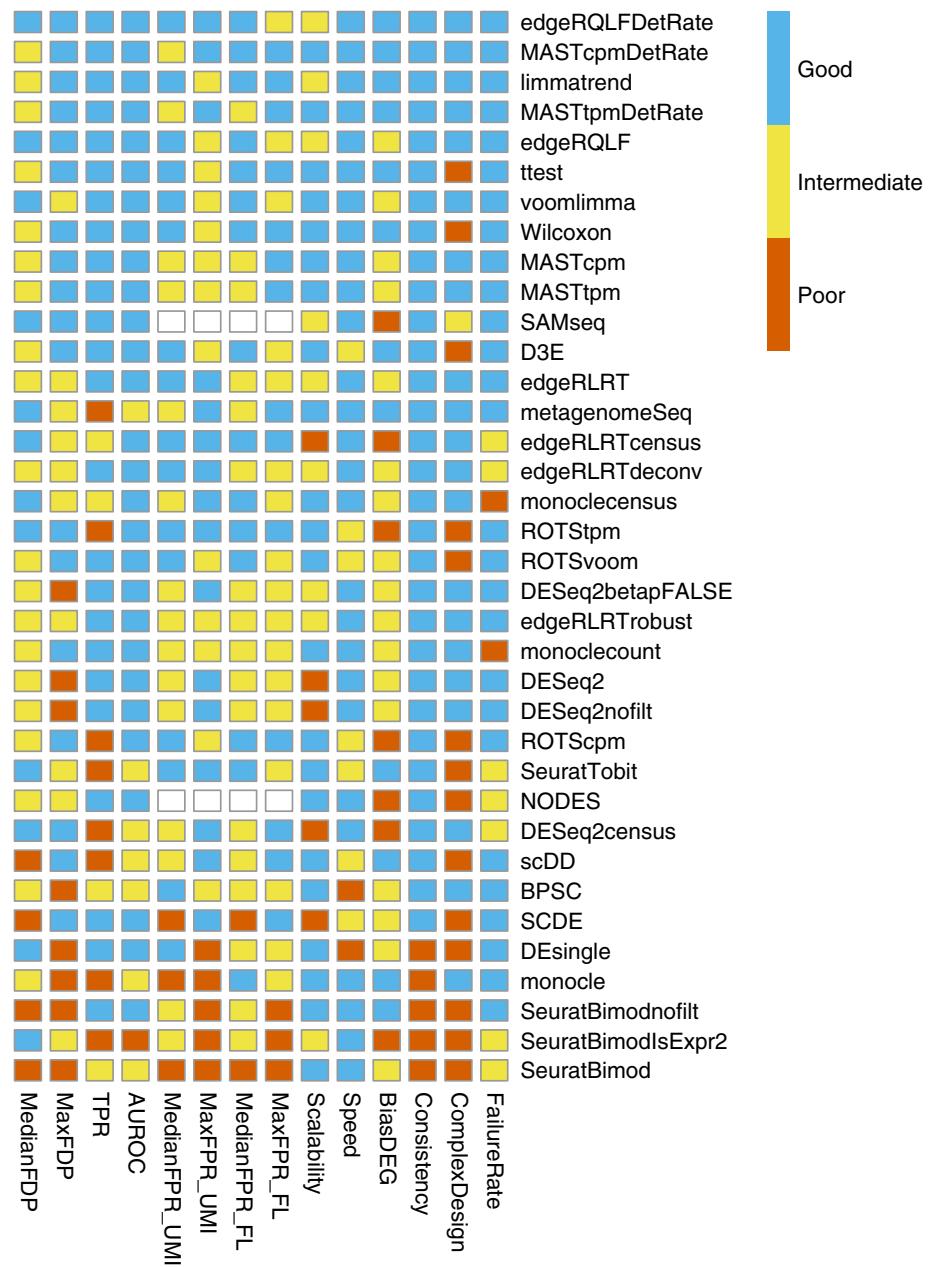
# The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells

Cole Trapnell<sup>1,2,6</sup>, Davide Cacchiarelli<sup>1-3,6</sup>, Jonna Grimsby<sup>2</sup>, Prapti Pokharel<sup>2</sup>, Shuqiang Li<sup>4</sup>, Michael Morse<sup>1,2</sup>, Niall J Lennon<sup>2</sup>, Kenneth J Livak<sup>4</sup>, Tarjei S Mikkelsen<sup>1-3</sup> & John L Rinn<sup>1,2,5</sup>

- Monocle: part of “tuxedo suite” for scRNA-seq analysis.
- Works for DE and clustering.
- Main idea for DE:
  - Model data with observed and dropout:  $Y = \begin{cases} Y^* & \text{if } Y^* > \lambda \\ \lambda & \text{if } Y^* \leq \lambda \end{cases}$
  - Use a generalized additive model (GAM) for design:
$$g(E(Y)) = \beta_0 + f_1(x_1) + f_2(x_2) + \cdots + f_m(x_m)$$
  - DE is tested from the GAM.

# scRNA-seq DE comparison

- A very informative comparison paper is Soneson and Robinson (2018) Nature Method.



# Cell clustering

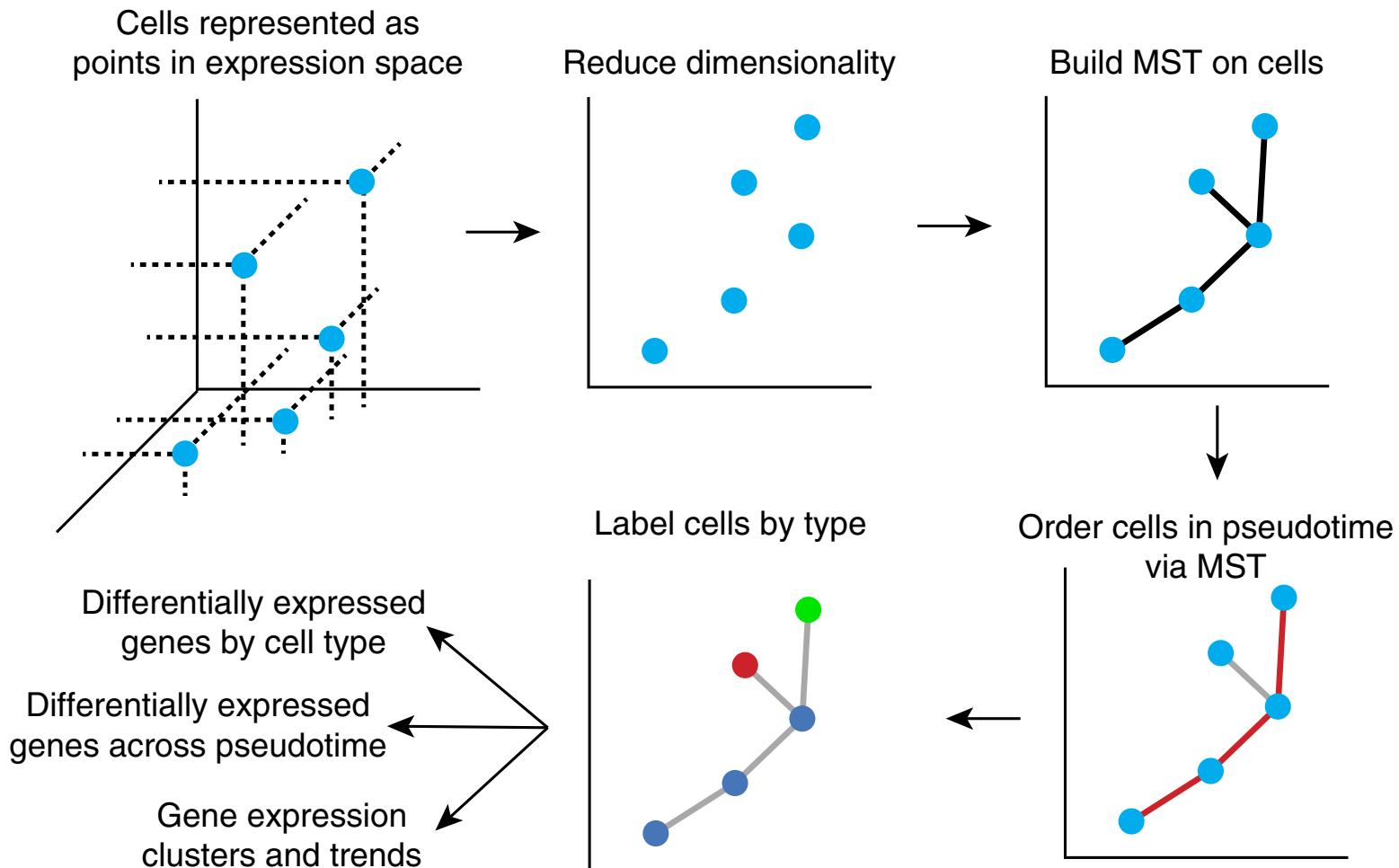
- The goals include:
  - Cluster cells into subgroups.
  - Model temporal transcriptomic dynamics: reconstruct “pseudo-time” for cells. This is useful for understanding development or disease progression.
- Traditional method like k-means or hierarchical clustering need to be used with caution due to dropout events.

# Clustering tools for scRNA-seq

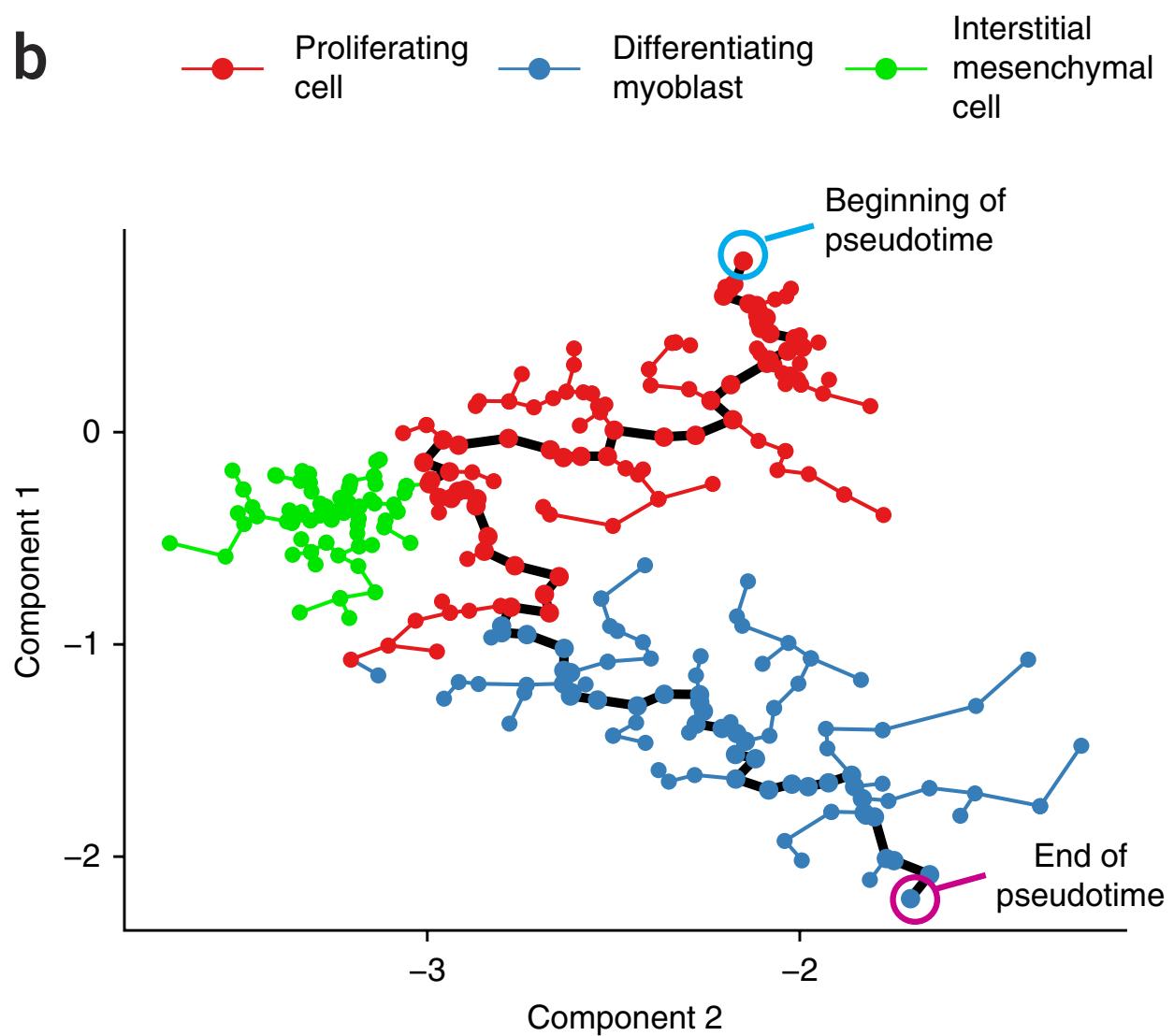
- A long list of tools: Seurat, TSCAN, SC3, CIDR, Monocle, ...
- Ideas are similar:
  - Select informative genes.
  - Dimension reduction of GE.
  - Cluster the cells based on reduced data.

# Monocle

a



# Monocle result



# Use Monocle Bioconductor package

First create a CellDataSet object using newCellDataSet function, then:

- Differential expression using `differentialGeneTest`.
- Cell ordering (pseudo-time estimation). This contains three steps:
  - Select a list of genes (often the DE genes) used for cell ordering. Use `setOrderingFilter` function to set that.
  - Dimension reduction using `reduceDimension` function.
  - Cell ordering using `orderCells` function.

```
### Create data object
pd <- new("AnnotatedDataFrame", data = sample_sheet)
fd <- new("AnnotatedDataFrame", data = gene_annotation)
dataobj <- newCellDataSet(as.matrix(expr_matrix),
                         phenoData = pd, featureData = fd)

### DE test
diff_test_res <- differentialGeneTest(dataobj,
                                         fullModelFormulaStr=GE~cond",
                                         reducedModelFormulaStr="GE~1")

### cell ordering
ordering_genes <- row.names(subset(diff_test_res, qval < 0.1))
dataobj <- setOrderingFilter(dataobj, ordering_genes)
dataobj <- reduceDimension(dataobj, use_irlba=FALSE)
dataobj <- orderCells(dataobj, num_paths=2, reverse=TRUE)
plot_spanning_tree(dataobj)
```

# Use SC3

```
require(SC3)
sce = SingleCellExperiment(
    assays = list( counts = as.matrix(counts) ,
                  logcounts = log2(as.matrix(counts) + 1) )
)
rowData(sce)$feature_symbol <- rownames(sce)
sce = sce[!duplicated(rowData(sce)$feature_symbol), ]
sce = sc3_prepare(sce)
if( missing(K) ) { ## estimate number of clusters if not given
    sce = sc3_estimate_k(sce)
    K = metadata(sce)$sc3$k_estimation
}
sce = sc3_calc_dists(sce)
sce = sc3_calc_transfs(sce)
sce = sc3_kmeans(sce, ks = K)
sce = sc3_calc_consens(sce)
colTb = colData(sce) [,1]
```

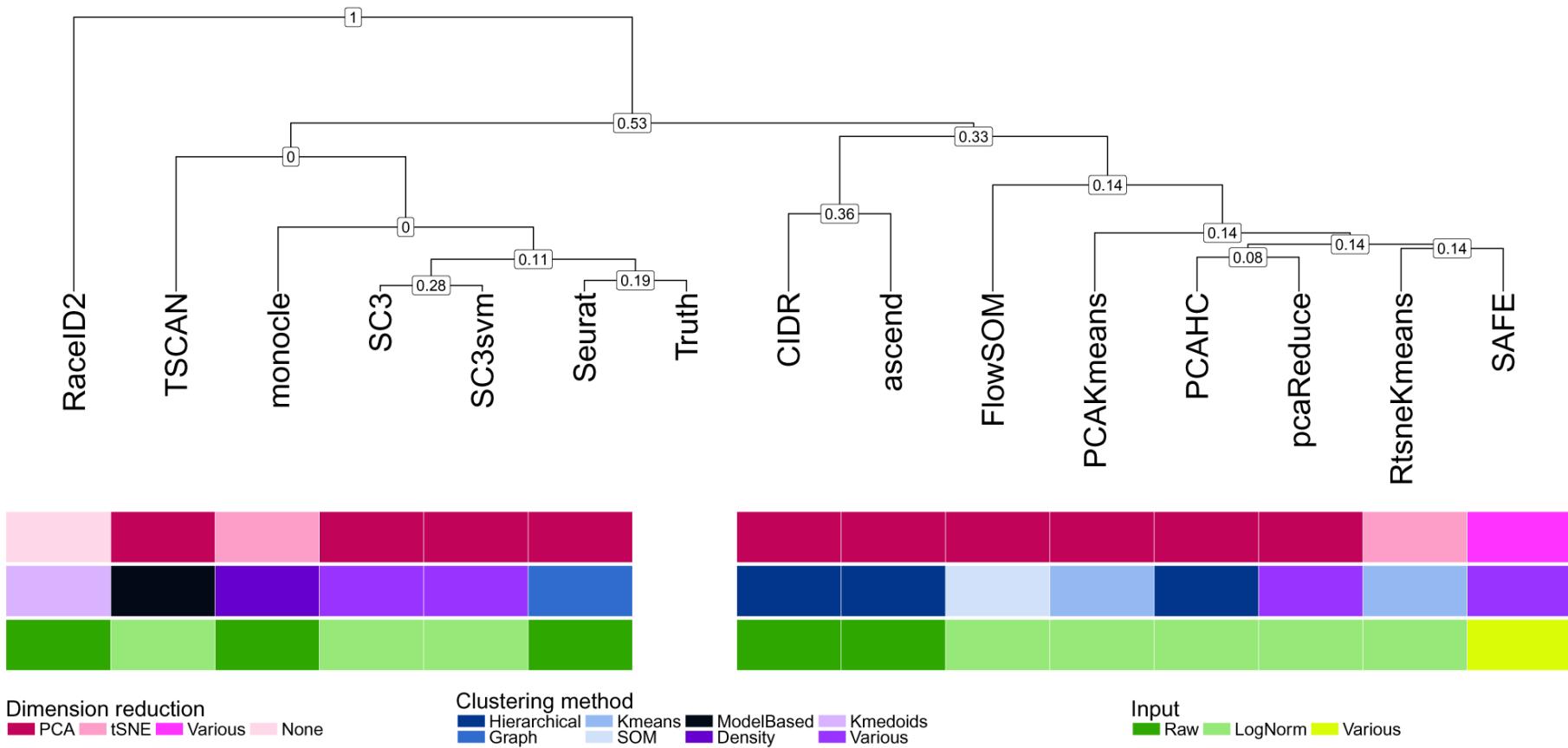
# Use Seurat

```
require(Seurat)

seuset = CreateSeuratObject( counts )
seuset = NormalizeData(object = seuset)
seuset = FindVariableFeatures(object = seuset)
seuset = ScaleData(object = seuset)
seuset = RunPCA(object = seuset)
seuset = FindNeighbors(object = seuset)
seuset = FindClusters(object = seuset)
return(seuset@active.ident)
```

# Clustering method comparison

- Duo et al. (2018) F1000



# **Detect rare cell type**

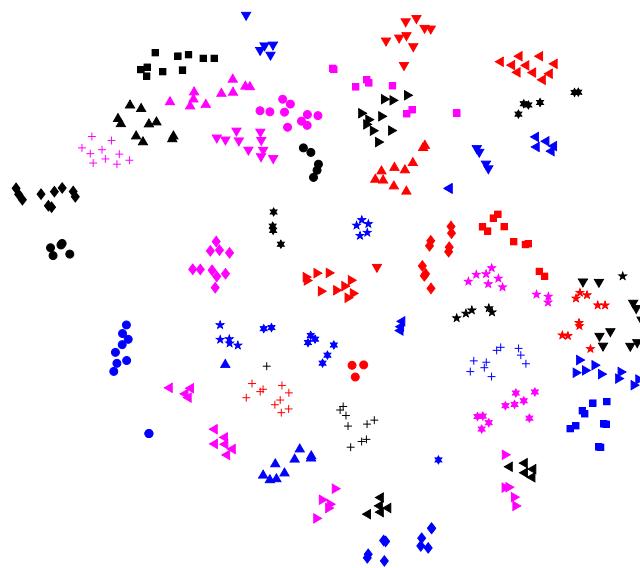
- Rare cell are “outliers” in the data.
- RacelD (Grun et al. 2015 Nature):
  - Normalize and log-transformed data.
  - Filter cells and genes
  - K-means clustering
  - Detect outliers from k-means result.
- GiniClust (Jiang et al. 2016 GB):
  - Difference is the gene filtering. It uses gini-index instead of variance to select genes.

# Missing data imputation

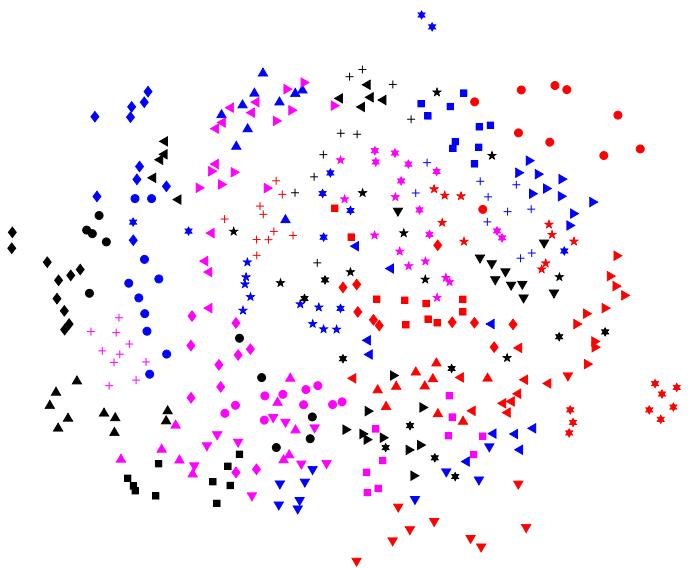
- Many zeros in scRNA-seq data.
- Attempts to impute (fill in) the zeros, with information from similar cells/genes.
- Available tools: sclImpute, drimpute, PBLR, SAVER, McImpute, netSmooth, ...

# t-SNE: a useful visualization tool

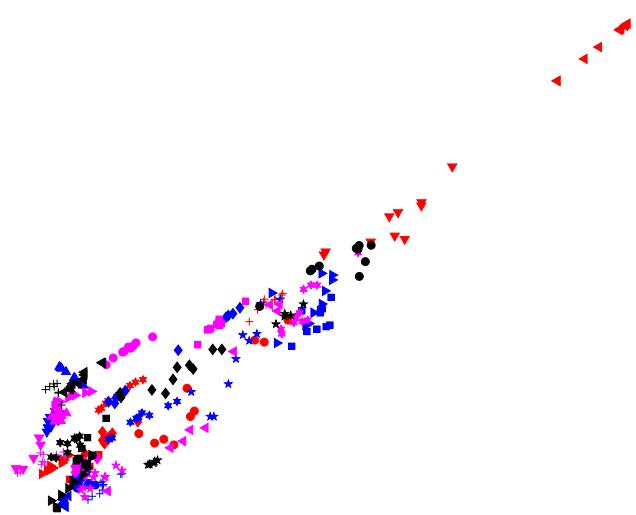
- t-SNE (t-distributed stochastic neighbor embedding): visualize high-dimensional data on 2-/3-D map.
- When project high-dimensional data into lower dimensional space, preserve the distances among data points.
  - This alleviate the problem that many clusters overlap on low dimensional space.
- Try to make the pairwise distances of points similar in high and low dimension.
- This is used in almost all scRNA-seq data visualization.
- Has “tsne” package on CRAN.



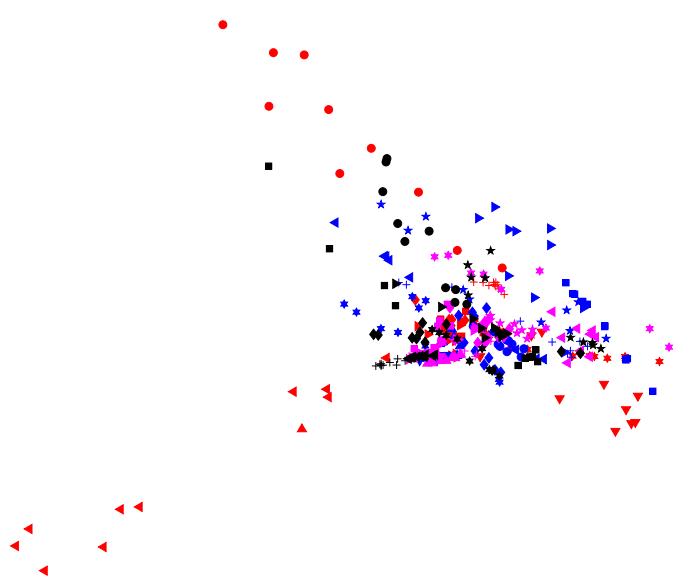
(a) Visualization by t-SNE.



(b) Visualization by Sammon mapping.



(c) Visualization by Isomap.



(d) Visualization by LLE.

# Summary for scRNA-seq

- The main interests are inter-cellular heterogeneity, expression dynamics, cell type discovery.
- Statistical questions include normalization, differential expression and clustering.
- Rooms for model development.

# Joint profiling in single cell

SINGLE-CELL GENOMICS

## Joint profiling of chromatin accessibility and gene expression in thousands of single cells

Junyue Cao<sup>1,2</sup>, Darren A. Cusanovich<sup>1\*†</sup>, Vijay Ramani<sup>1\*</sup>, Delasa Aghamirzaie<sup>1</sup>, Hannah A. Pliner<sup>1</sup>, Andrew J. Hill<sup>1</sup>, Riza M. Daza<sup>1</sup>, Jose L. McFaline-Figueroa<sup>1</sup>, Jonathan S. Packer<sup>1</sup>, Lena Christiansen<sup>3</sup>, Frank J. Steemers<sup>3</sup>, Andrew C. Adey<sup>4,5</sup>, Cole Trapnell<sup>1,6,7‡</sup>, Jay Shendure<sup>1,6,7,8‡</sup>

ARTICLE

DOI: 10.1038/s41467-018-03149-4

OPEN

## scNMT-seq enables joint profiling of chromatin accessibility DNA methylation and transcription in single cells

Stephen J. Clark  <sup>1</sup>, Ricard Argelaguet<sup>2,3</sup>, Chantriolnt-Andreas Kapourani  <sup>4</sup> Thomas M. Stubbs<sup>1</sup>, Heather J. Lee<sup>1,5,6</sup>, Celia Alda-Catalinas  <sup>1</sup>, Felix Krueger  <sup>7</sup> Guido Sanguinetti<sup>4</sup>, Gavin Kelsey  <sup>1,8</sup> John C. Marioni  <sup>2,3,5</sup> Oliver Stegle  <sup>2</sup> Wolf Reik<sup>1,5,8</sup>

# Single cell GE microarray



## Single cell-derived clonal analysis of human glioblastoma links functional and genomic heterogeneity

Mona Meyer<sup>a,1</sup>, Jüri Reimand<sup>b,c,1</sup>, Xiaoyang Lan<sup>a,b</sup>, Renee Head<sup>a</sup>, Xueming Zhu<sup>a</sup>, Michelle Kushida<sup>a</sup>, Jessica C. Pressey<sup>e</sup>, Anath C. Lionel<sup>b,f</sup>, Ian D. Clarke<sup>a,g</sup>, Michael Cusimano<sup>h</sup>, Jeremy A. Squire<sup>i</sup>, Stephen Bernstein<sup>j</sup>, Melanie A. Woodin<sup>e</sup>, Gary D. Bader<sup>b,c,2</sup>, and Peter B. Dirks<sup>a,b,k,2</sup>

<sup>a</sup>Division of Neurosurgery, Program in Developmental and Stem Cell Biology, Arthur and Sonia Labatt Brain Tumour Research Ce

# Single cell lncRNA

**Genome Biology**



This Provisional PDF corresponds to the article as it appeared upon acceptance. Fully formatted PDF and full text (HTML) versions will be made available soon.

## **Localization and abundance analysis of human lncRNAs at single-cell and single-molecule resolution**

*Genome Biology*

doi:10.1186/s13059-015-0586-4

Moran N Cabili ([nmcabili@broadinstitute.org](mailto:nmcabili@broadinstitute.org))  
Margaret C Dunagin ([dunagin@seas.upenn.edu](mailto:dunagin@seas.upenn.edu))  
Patrick D McClanahan ([pmcll@seas.upenn.edu](mailto:pmcll@seas.upenn.edu))  
Andrew Biaesch ([biaesch@gmail.com](mailto:biaesch@gmail.com))  
Olivia Padovan-Merhar ([opadovan@sas.upenn.edu](mailto:opadovan@sas.upenn.edu))  
Aviv Regev ([aregev@broadinstitute.org](mailto:aregev@broadinstitute.org))  
John L Rinn ([john\\_rinn@harvard.edu](mailto:john_rinn@harvard.edu))  
Arjun Raj ([arjunraj@seas.upenn.edu](mailto:arjunraj@seas.upenn.edu))

Published online: 29 January 2015

# Grand summary for scSeq

- Single-cell biology reveals a lot of information that can't be detected from bulk data.
- Data are much noisier, and more difficult to analyze.
- Some room for method development, but very competitive.