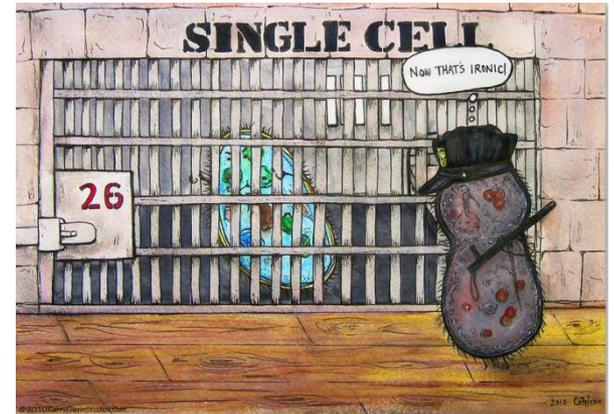


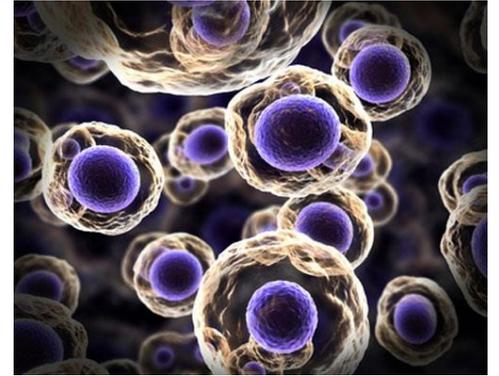
# Single cell transcriptomics

Ester Feldmesser

Introduction to Deep Sequencing Analysis—  
June 2015



# Single cell applications



- ▶ Single-cell genome sequencing and structure
- ▶ Single-cell methylation analysis
- ▶ Single-cell transcriptomics (and proteomics)
- ▶ Combined approaches

# Overview

- ▶ Motivation for single cell transcriptomics
  - ▶ Protocols
  - ▶ Pitfalls and how to overcome them
  - ▶ One specific protocol: Quantitative single-cell RNA-seq with unique molecular identifiers
  - ▶ Some examples of results
  - ▶ *In situ* single cell transcriptomics
- 



# Reference papers

## REVIEWS

### SINGLE-CELL OMICS

## Computational and analytical challenges in single-cell transcriptomics

*Oliver Stegle<sup>1</sup>, Sarah A. Teichmann<sup>1,2</sup> and John C. Marioni<sup>1,2</sup>*

**Abstract** | The development of high-throughput RNA sequencing (RNA-seq) at the single-cell level has already led to profound new discoveries in biology, ranging from the identification of novel cell types to the study of global patterns of stochastic gene expression. Alongside the technological breakthroughs that have facilitated the large-scale generation of single-cell transcriptomic data, it is important to consider the specific computational and analytical challenges that still have to be overcome. Although some tools for analysing RNA-seq data from bulk cell populations can be readily applied to single-cell RNA-seq data, many new computational strategies are required to fully exploit this data type and to enable a comprehensive yet detailed study of gene expression at the single-cell level.

doi:10.1038/nrg3833

## Quantitative single-cell RNA-seq with unique molecular identifiers

Saiful Islam, Amit Zeisel, Simon Joost, Gioele La Manno, Pawel Zajac, Maria Kasper, Peter Lönnerberg & Sten Linnarsson

[Affiliations](#) | [Contributions](#) | [Corresponding author](#)

*Nature Methods* **11**, 163–166 (2014) | doi:10.1038/nmeth.2772

Received 27 September 2013 | Accepted 25 November 2013 | Published online 22 December 2013



Single-cell RNA sequencing (RNA-seq) is a powerful tool to reveal cellular heterogeneity, discover new cell types and characterize tumor microevolution. However, losses in cDNA synthesis and bias in cDNA amplification lead to severe quantitative errors. We show that molecular labels—random sequences that label individual molecules—can nearly

nized

### APPLICATIONS OF NEXT-GENERATION SEQUENCING

## Single-cell sequencing-based technologies will revolutionize whole-organism science

*Ehud Shapiro<sup>1,2</sup>, Tamir Biezuner<sup>1,2</sup> and Sten Linnarsson<sup>3</sup>*

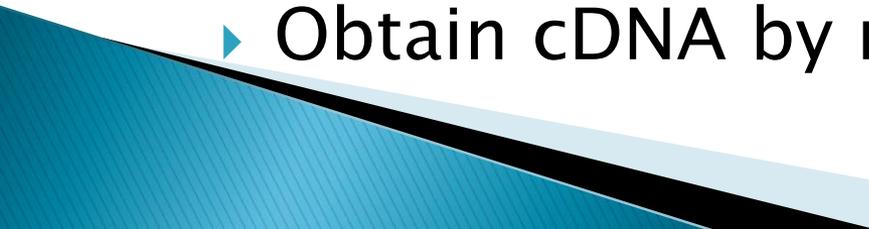
**Abstract** | The unabated progress in next-generation sequencing technologies is fostering a wave of new genomics, epigenomics, transcriptomics and proteomics technologies. These sequencing-based technologies are increasingly being targeted to individual cells, which will allow many new and longstanding questions to be addressed. For example, single-cell genomics will help to uncover cell lineage relationships; single-cell transcriptomics will supplant the coarse notion of marker-based cell types; and single-cell epigenomics and proteomics will allow the functional states of individual cells to be analysed. These technologies will become integrated within a decade or so, enabling high-throughput, multi-dimensional analyses of individual cells that will produce detailed knowledge of the cell lineage trees of higher organisms, including humans. Such studies will have important implications for both basic biological research and medicine.

[www.nature.com/reviews/genetics](http://www.nature.com/reviews/genetics)

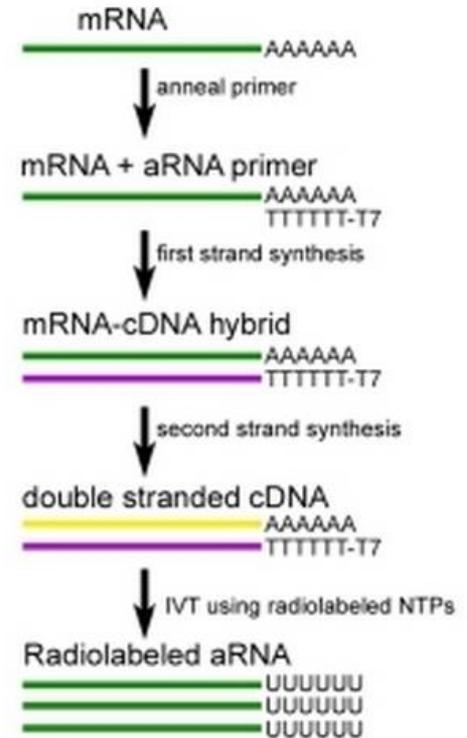
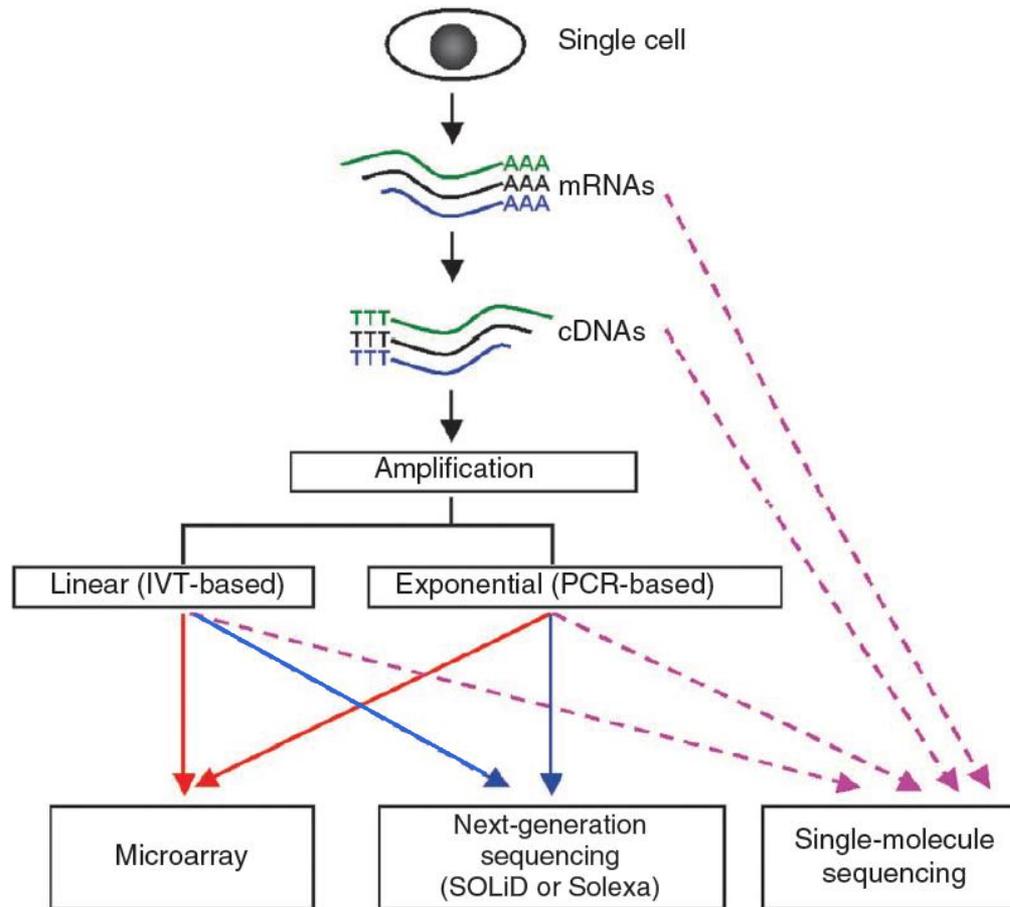
# Motivation

- ▶ Tissues are complex and not homogenous
  - ▶ Different cell types gene expression can be masked
  - ▶ Identification of cell types and cell states
  - ▶ Learn about noise in gene expression
  - ▶ Ideally would like to combine transcriptional profiling with spatial information
- 

# Protocols for single cell

- ▶ Isolate individual cells:
    - Manually by micromanipulation
    - Manually by fluorescence-activated cell sorting (FACS)
    - Exploiting a microfluidics based system
    - Laser capture microdissection
  - ▶ Capture the poly-adenylated fraction of mRNA molecules
  - ▶ Obtain cDNA by reverse transcription
- 

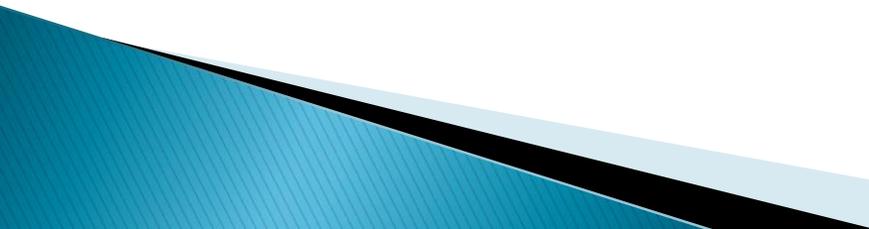
# Single-cell transcriptome



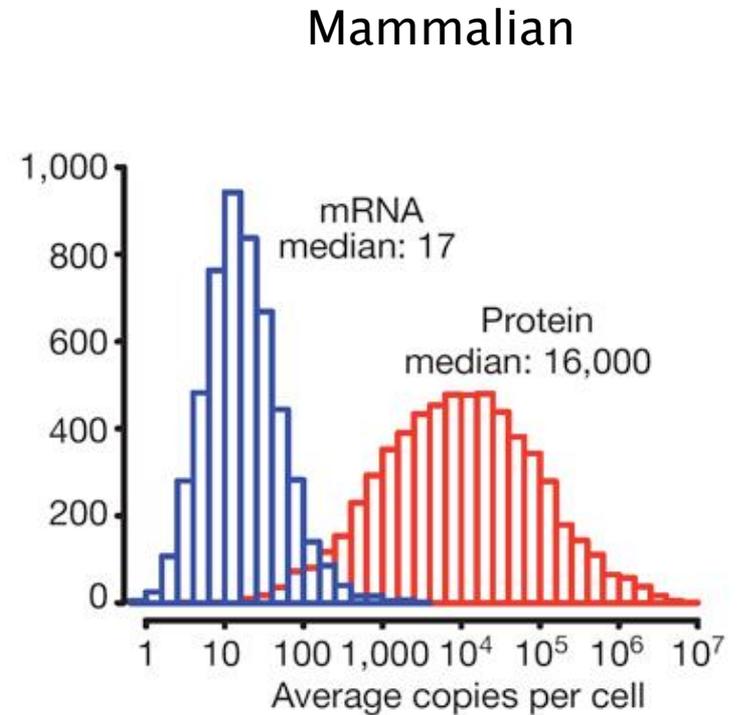
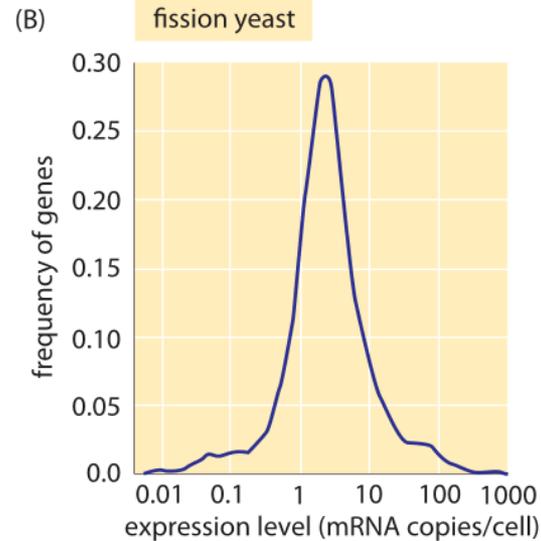
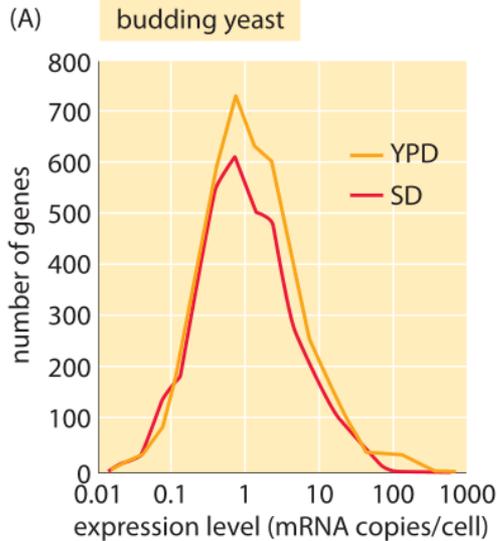
# Overview

- ▶ Motivation for single cell transcriptomics
  - ▶ Protocols
  - ▶ Pitfalls and how to overcome them
  - ▶ One specific protocol: Quantitative single-cell RNA-seq with unique molecular identifiers
  - ▶ Some examples of results
  - ▶ In situ single cell transcriptomics
- 

# How many mRNAs are there in a cell?

- ▶ Approximately 360,000 mRNA molecules in a single mammalian cell, made up of approximately 12,000 different transcripts
  - ▶ Some mRNAs comprise 3% of the mRNA pool whereas others account for less than 0.1%. These rare or low-abundance mRNAs may have a copy number of only 5–15 molecules per cell.
- 

# mRNAs numbers in a cell



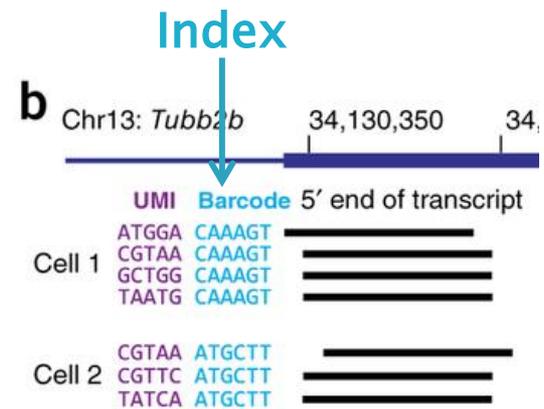
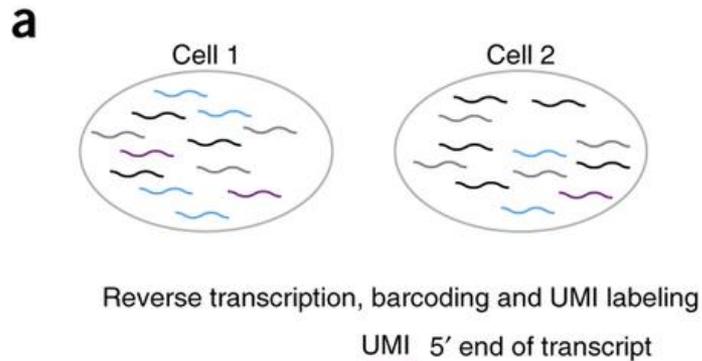
# Incorporating quantitative standards

- ▶ Use extrinsic spike-in molecules: artificial spike-in mix is the External RNA Control Consortium (ERCC) set of 92 synthetic spikes based on bacterial sequences
- ▶ Unique molecular identifiers (UMIs) have been used to barcode individual molecules, when sequencing the 3' or 5' end of the amplified transcript.

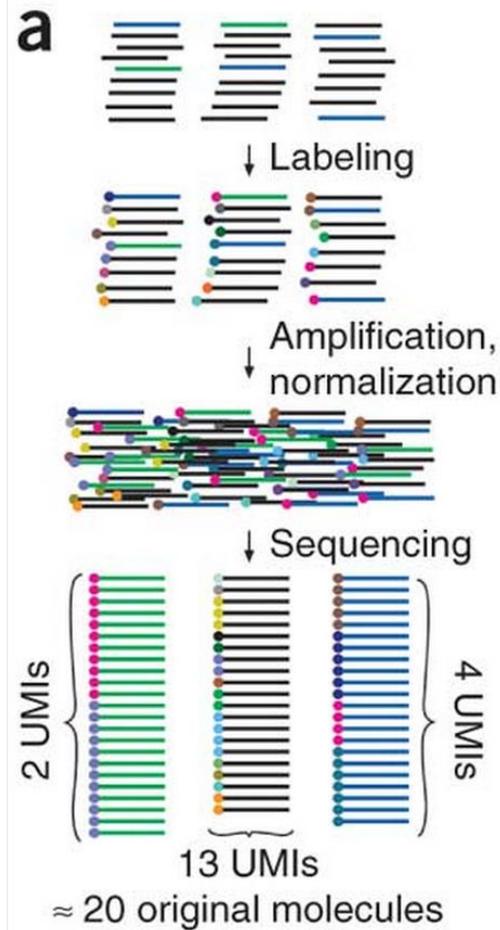
# Unique molecular identifier (UMI)



- Universal Adapter
- DNA Fragment of Interest
- Indexed Adapter
- 6 Base Index Region
- UMI



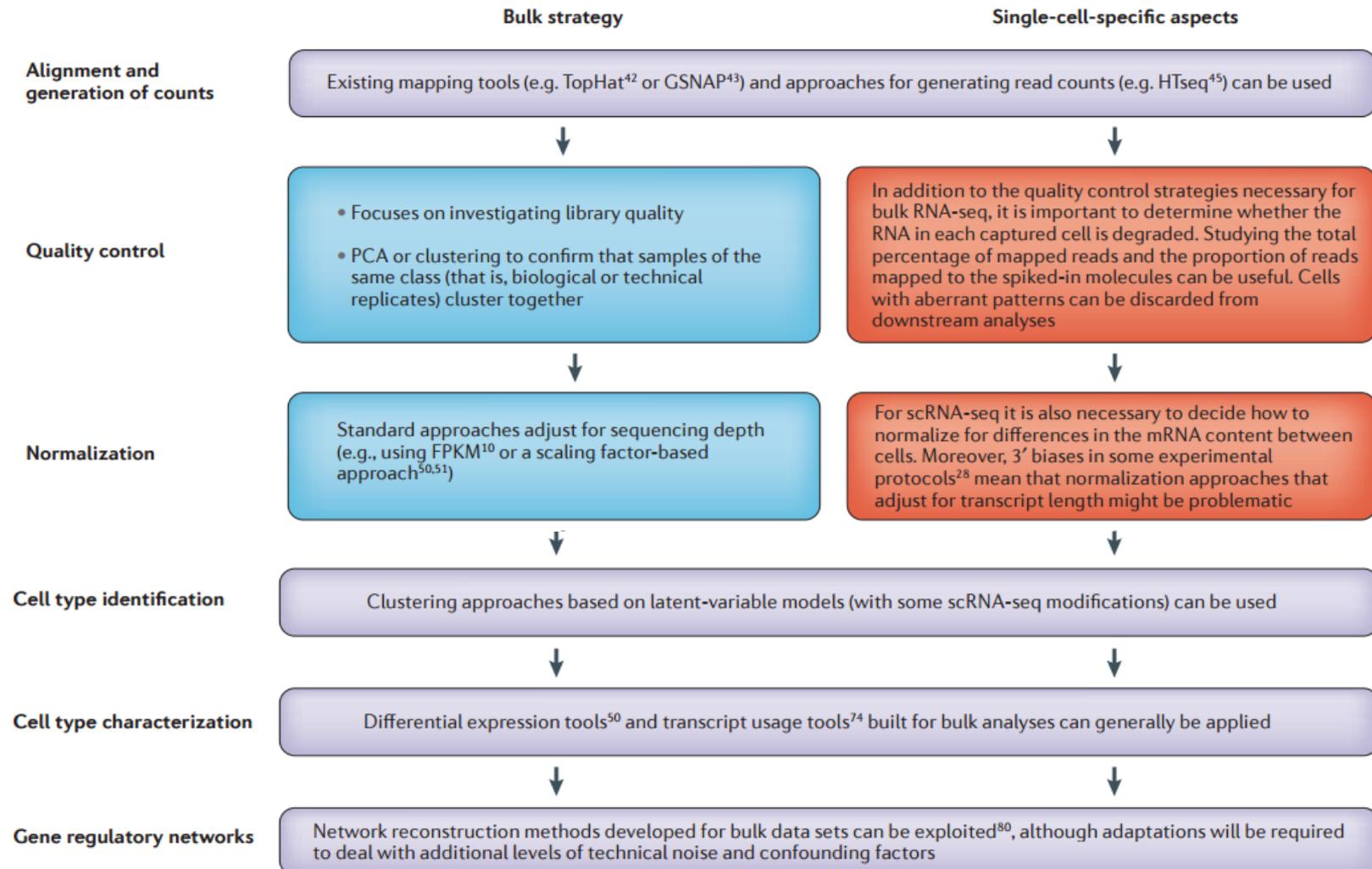
# UMIs can be generated by adding oligonucleotide labels



# Transcript quantification and quality control



# Comparison of bulk and scRNA-seq analytical strategies

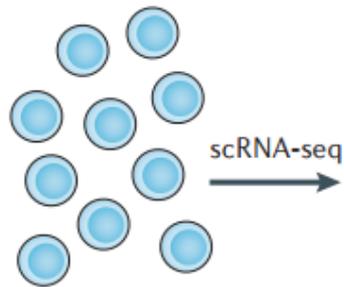


# Quality control

- ▶ Cells captured may contain **degraded RNA**: for example, because the cell is **stressed** because of extraction and isolation
- ▶ Check:
  - ✓ Mapping to the genome
  - ✓ Ratio of reads mapped to the genome versus the number of reads mapped to the extrinsic spike-ins
  - ✓ Principal component analysis
- ▶ **Discard outlier cells**

# Quality control (2)

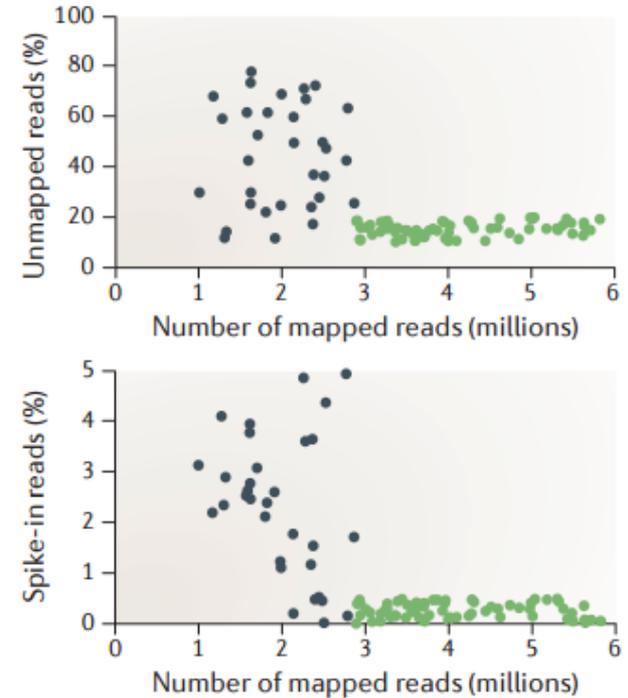
a



Read counts

	Cell 1	Cell 2	...
Gene 1	25	918	
Gene 2	0	456	
...			
Spike 1	103	180	
Spike 2	13	19	
...			

- Poor-quality cells (high percentage of unmapped reads or spike-in reads)
- Higher-quality cells



## b. PCA

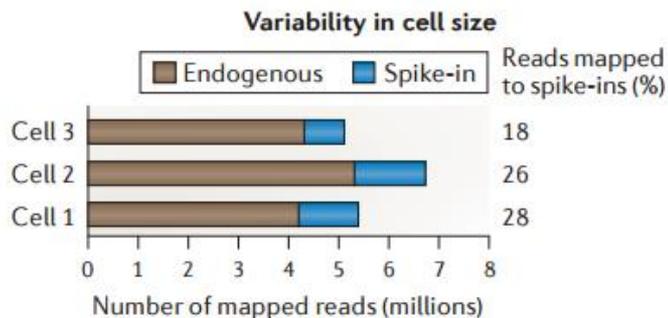
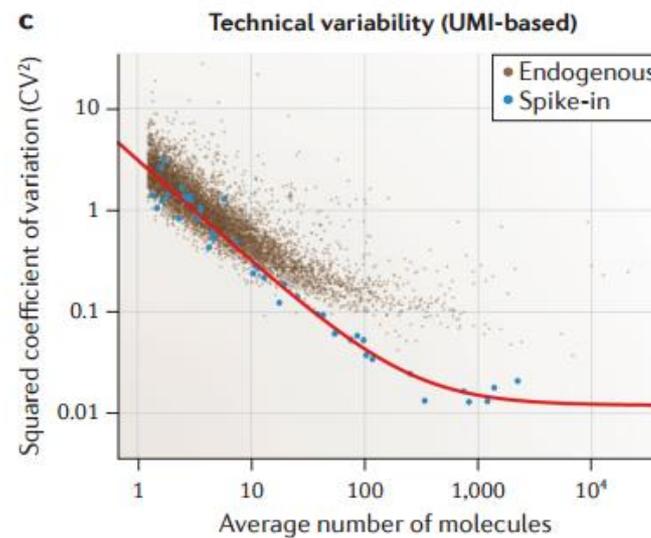
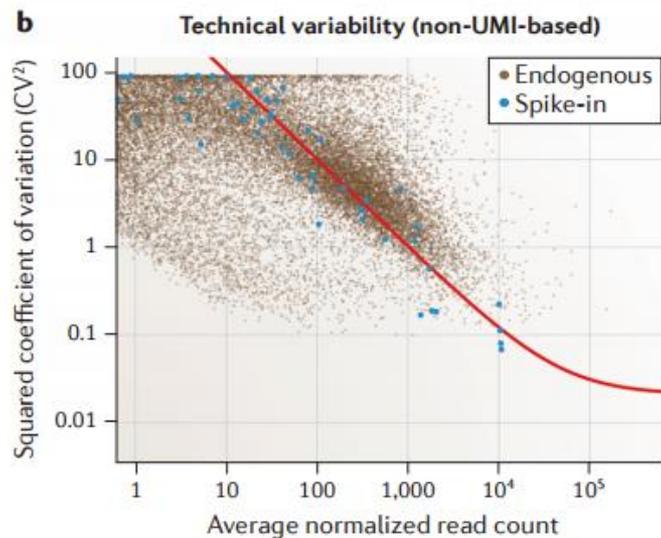
The expectation when applying PCA is that good-quality cells cluster together and poor-quality cells are outliers.

Poor-quality cells may also form a second distinct cluster.

# Normalization method

- ▶ Common approaches for normalizing bulk RNA-seq data make an implicit assumption: that the total amount of RNA processed in each sample is approximately the same or that the variation is technical.
- ▶ This assumption has been shown to be misleading (for example, upregulation of *MYC* leads to a two-fold increase in the number of transcripts)

# Normalization



**Capture efficiency**

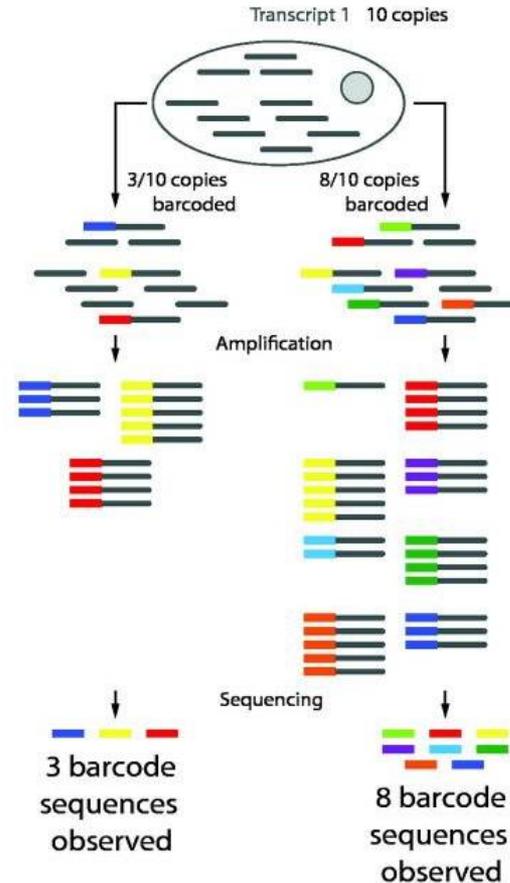
Cell	Number of UMIs associated with spike-in molecules	Initial number of spiked in molecules	Capture efficiency (%)
1	1,436	10,000	14.36
2	2,598	10,000	25.98
3	987	10,000	9.87

# Normalization issues

- ▶ ERCC have short poly-A tails (efficiency RT)
- ▶ ERCC spikes are short (500 – 2000)
- ▶ 5'-to-3' length bias is inherent to many scRNA-seq protocols (problems normalizing)
- ▶ Normalizing for transcript length is challenging with current scRNA-seq protocols (3' biases)
- ▶ Differences in the efficiency of the reverse transcription reaction (UMI)

# Differences in the efficiency of the reverse transcription reaction (UMI)

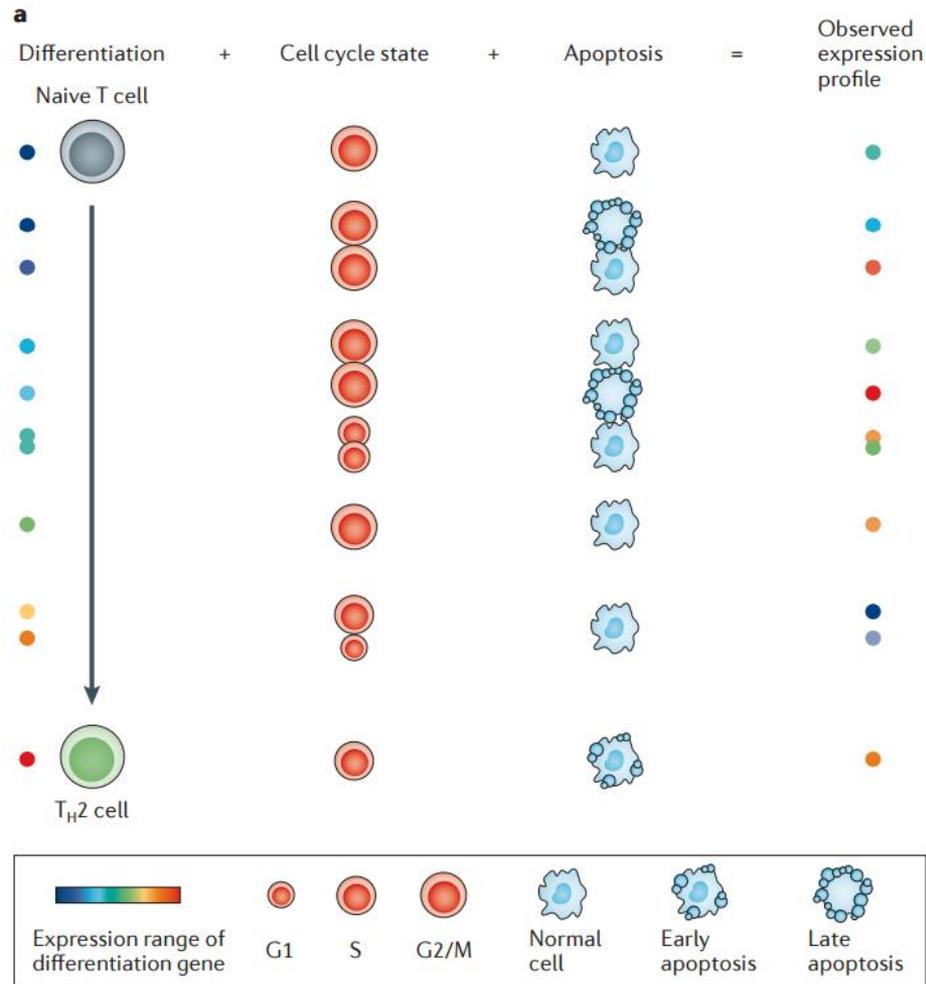
UMI principle and the output depending on reaction efficiency



# scRNA-seq experimental design

- ▶ Number of cells sequenced
- ▶ Depth (# reads) for each individual cell
  - Jaitin et al generated an average of 22,000 aligned sequence reads from 1,536 cells
  - Mahata et al. generated 12–20 million sequence reads from each of 93 cells
- ▶ Batch effects
- ▶ Biological factors: cell cycle, physical condition of the cell

# Confounding variables



# Overview

- ▶ Motivation for single cell transcriptomics
  - ▶ Protocols
  - ▶ Pitfalls and how to overcome them
  - ▶ One specific protocol: Quantitative single-cell RNA-seq with unique molecular identifiers
  - ▶ Some examples of results
  - ▶ *In situ* single cell transcriptomics
- 



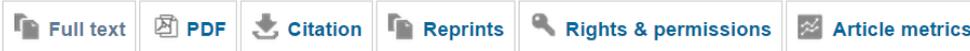
# Quantitative single-cell RNA-seq with unique molecular identifiers

Saiful Islam, Amit Zeisel, Simon Joost, Gioele La Manno, Pawel Zajac, Maria Kasper, Peter Lönnerberg & Sten Linnarsson

[Affiliations](#) | [Contributions](#) | [Corresponding author](#)

*Nature Methods* **11**, 163–166 (2014) | doi:10.1038/nmeth.2772

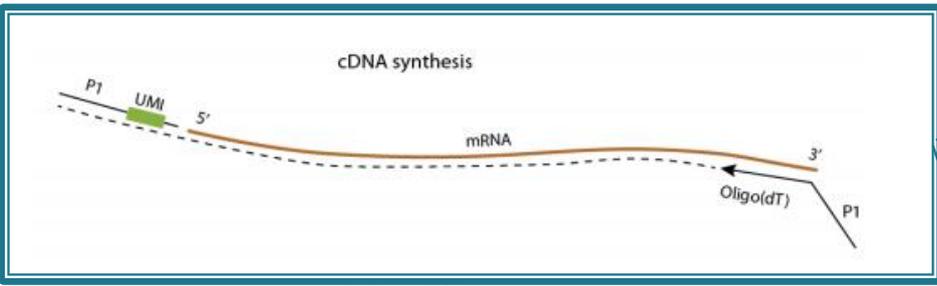
Received 27 September 2013 | Accepted 25 November 2013 | Published online 22 December 2013



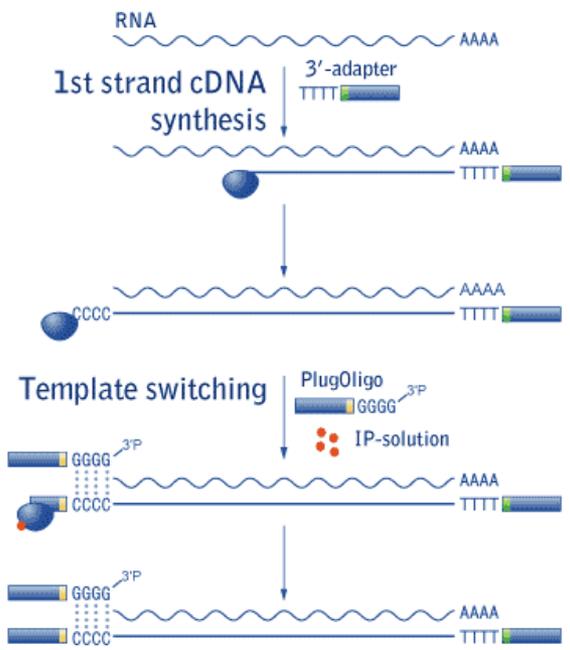
Single-cell RNA sequencing (RNA-seq) is a powerful tool to reveal cellular heterogeneity, discover new cell types and characterize tumor microevolution. However, losses in cDNA synthesis and bias in cDNA amplification lead to severe quantitative errors. We show that molecular labels—random sequences that label individual molecules—can nearly eliminate amplification noise, and that microfluidic sample preparation and optimized reagents produce a fivefold improvement in mRNA capture efficiency.

- ▶ Experiment performed on ES cells
- ▶ Unique molecular identifiers – reduce amplification noise
- ▶ Microfluidic sample preparation – improves mRNA capture

# Unique molecular identifier (UMI) method



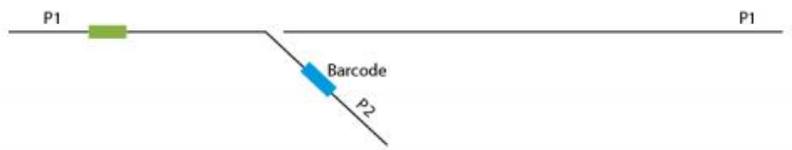
## Template switch



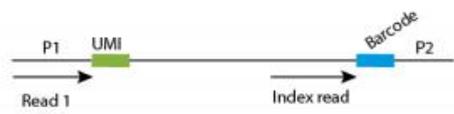
## PCR amplification



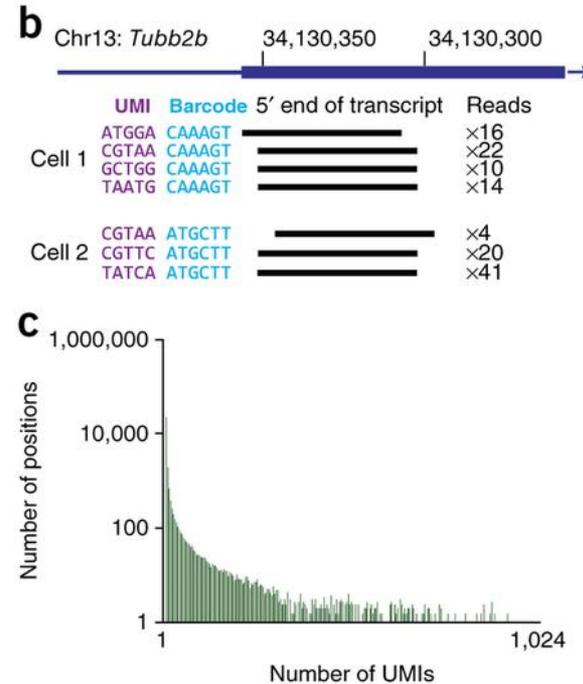
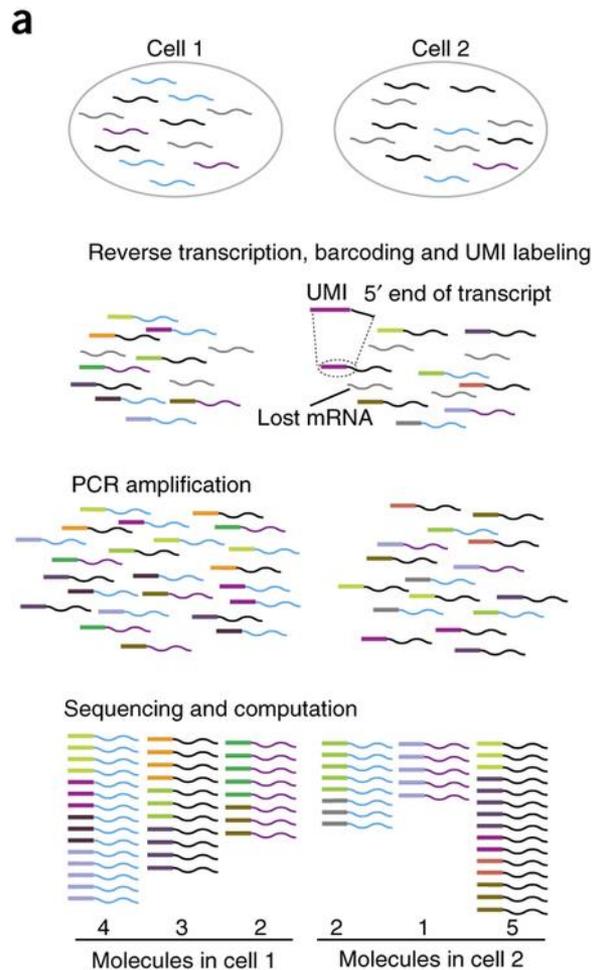
## Tagmentation



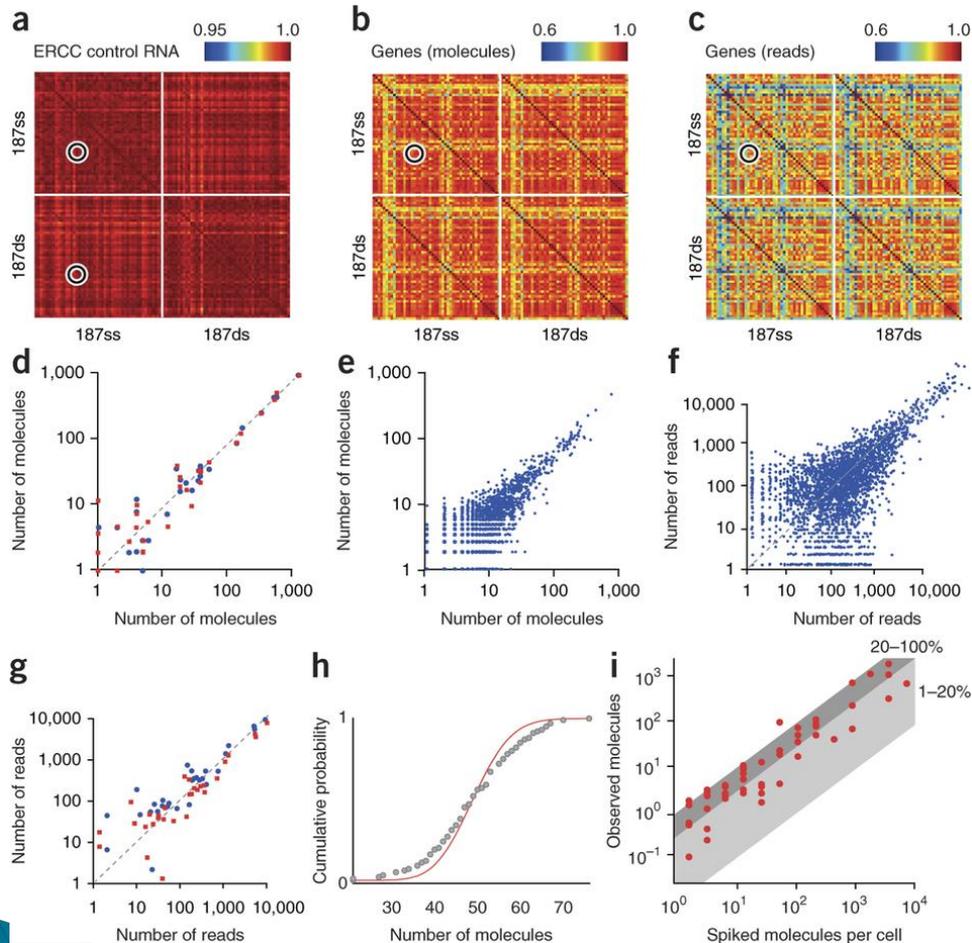
## Sequencing



# Overview of tagging single mRNA molecules



# Reproducibility of molecule counting



(a,b,c) Pairwise correlation coefficients calculated for ERCC, endogenous genes using molecule counts or counting reads instead of molecules

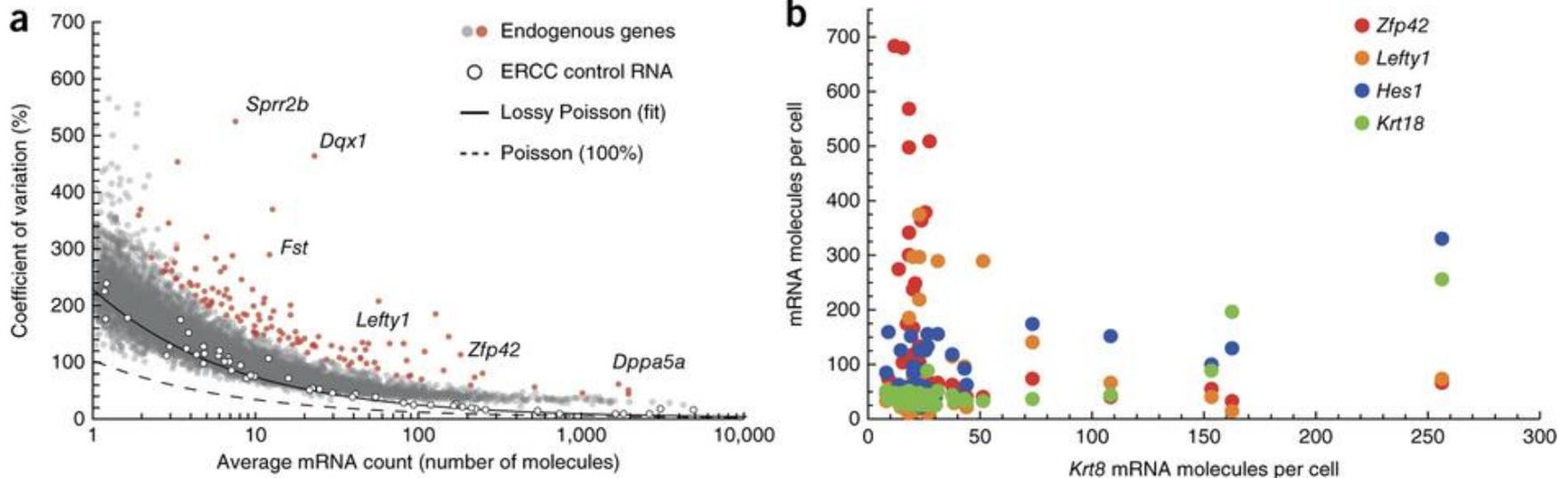
(d) Pairwise comparison of two wells indicated in a. Red squares and blue dots show comparisons within and between libraries, respectively.

(e,f) Scatterplot showing the two cells indicated in b or e, based on molecule counts or read counts respectively.

(g) Scatterplot as in d but using reads instead of molecules. (h) Distribution of molecule counts for a single ERCC spike-in transcript (gray dots) compared with the cumulative density function of the Poisson distribution (red line).

(i) mRNA capture efficiency shown as observed molecule counts versus number of spiked-in molecules for ERCC control RNA transcripts. The shaded bands indicate efficiencies above (dark gray) and below (light gray) 20%. Each red dot represents the average of a single ERCC RNA across 96 wells. Similar results were obtained in one replicate experiment.

# Transcriptional noise in ES cells



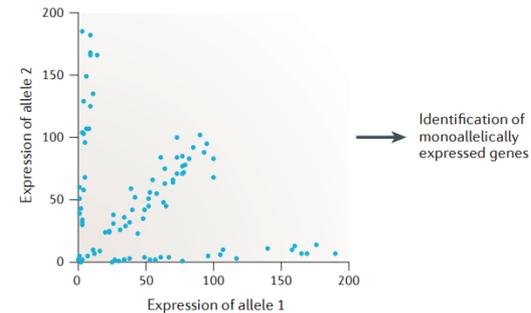
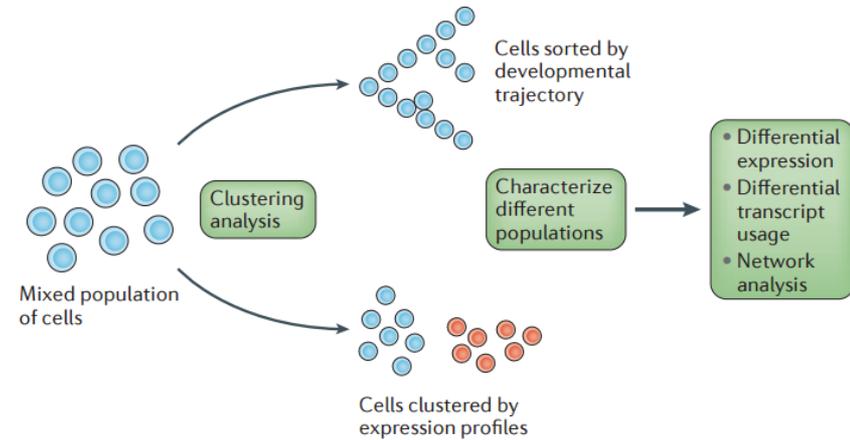
- (a) The coefficient of variation as a function of the mean number of mRNA molecules detected, for genes expressed in ES cells ( $n = 4$ ) Genes with significant ( $\alpha = 0.05$ ) excess noise are shown as red dots
- (b) Coexpression of noisy genes. The expression of four genes, across 41 ES cells, is plotted against the expression of *Krt8*. Two branches were observed, interpreted as 'epiblast-like' (high *Krt8*, *Krt18* and *Hes1*) and 'pluripotent-like' (high *Zfp42* and *Lefty1*)

# Overview

- ▶ Motivation for single cell transcriptomics
  - ▶ Protocols
  - ▶ Pitfalls and how to overcome them
  - ▶ One specific protocol: Quantitative single-cell RNA-seq with unique molecular identifiers
  - ▶ Some examples of results
  - ▶ *In situ* single cell transcriptomics
- 

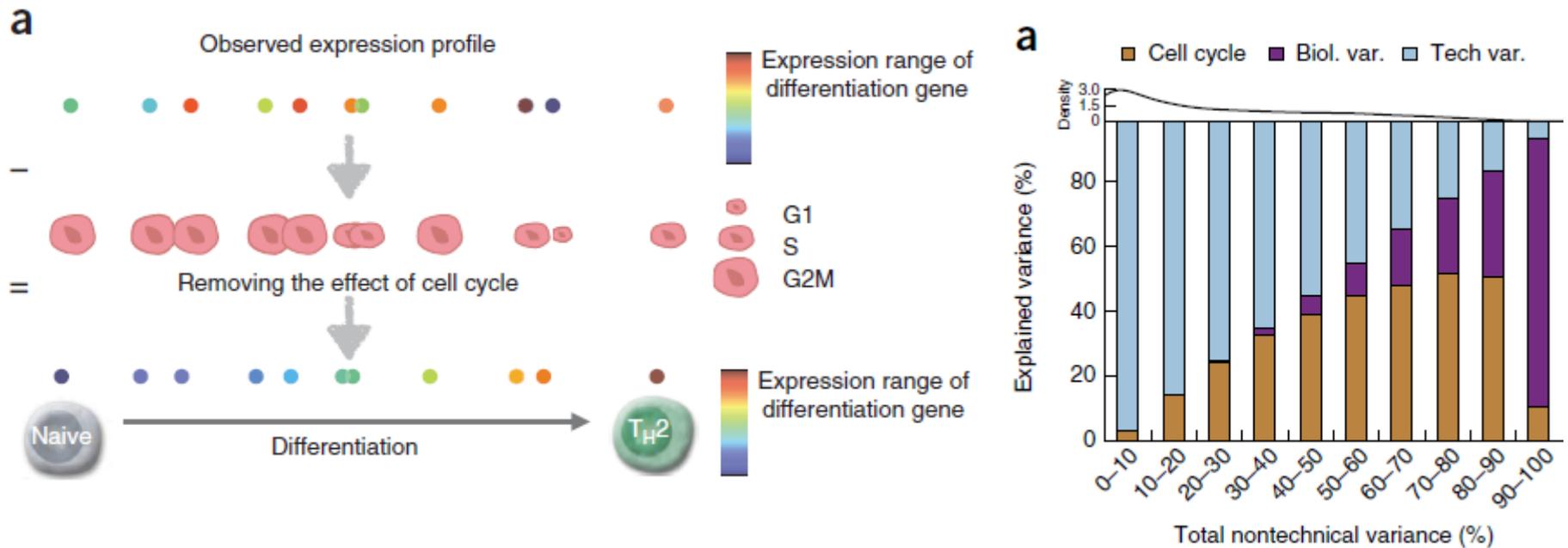
# What information can be obtained from scRNA-seq?

- ▶ Identification of cell type and cellular state
- ▶ Differential expression at the gene level
- ▶ Differential expression at the transcript-isoform level?
- ▶ Identification of highly variable genes
- ▶ Kinetics of gene expression
- ▶ Detection of allele-specific expression



# Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells

Florian Buettner<sup>1,2,5</sup>, Kedar N Natarajan<sup>2,3,5</sup>, F Paolo Casale<sup>2</sup>, Valentina Proserpio<sup>2,3</sup>, Antonio Scialdone<sup>2,3</sup>, Fabian J Theis<sup>1,4</sup>, Sarah A Teichmann<sup>2,3</sup>, John C Marioni<sup>2,3</sup> & Oliver Stegle<sup>2</sup>

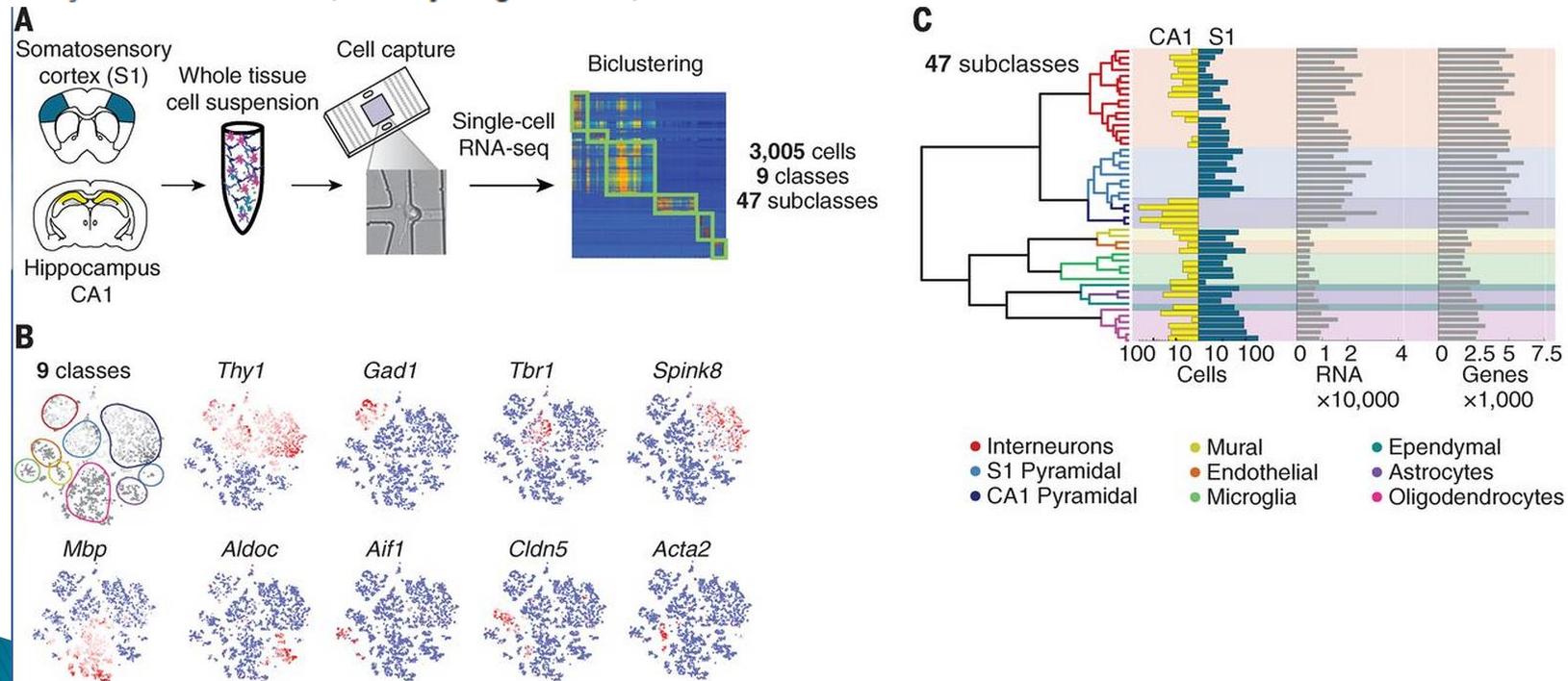


- Use a latent variable model approach to regress out variation attributable to the cell-cycle
- Technique can be used to remove variation for any given set or sets of genes

REPORT

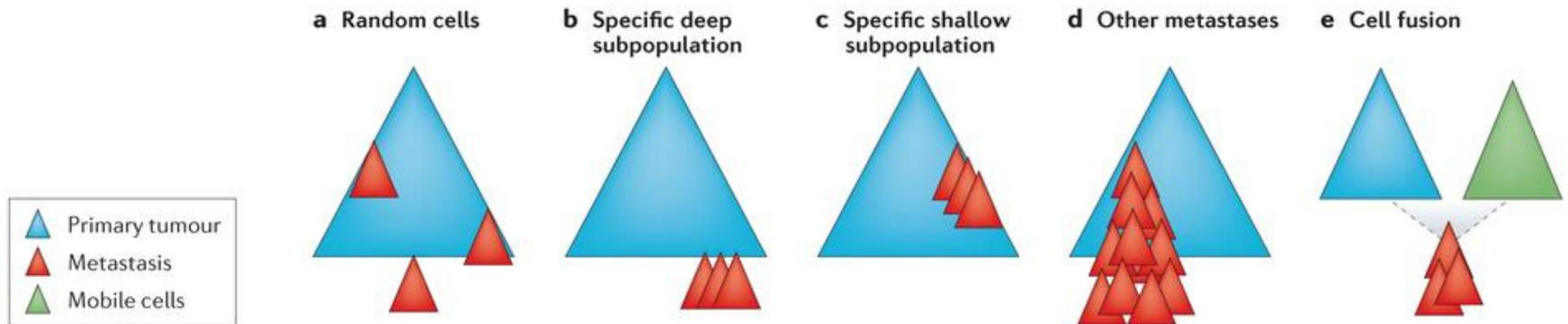
# Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq

Amit Zeisel<sup>1,\*</sup>, Ana B. Muñoz-Manchado<sup>1,\*</sup>, Simone Codeluppi<sup>1</sup>, Peter Lönnerberg<sup>1</sup>, Gioele La Manno<sup>1</sup>,  
Anna Juréus<sup>1</sup>, Sueli Marques<sup>1</sup>, Hermany Munguba<sup>1</sup>, Liqun He<sup>2</sup>, Christer Betsholtz<sup>2,3</sup>, Charlotte Rolny<sup>4</sup>,  
Gonçalo Castelo-Branco<sup>1</sup>, Jens Hjerling-Leffler<sup>1,†</sup>, Sten Linnarsson<sup>1,†</sup>



# Clinical use of single cell genomics/transcriptomics

*Cell lineage reconstruction of cancer will elucidate its development*

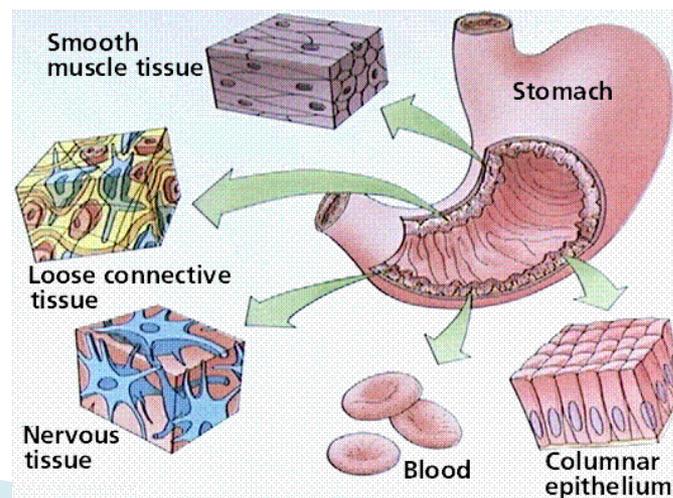


Nature Reviews | Genetics

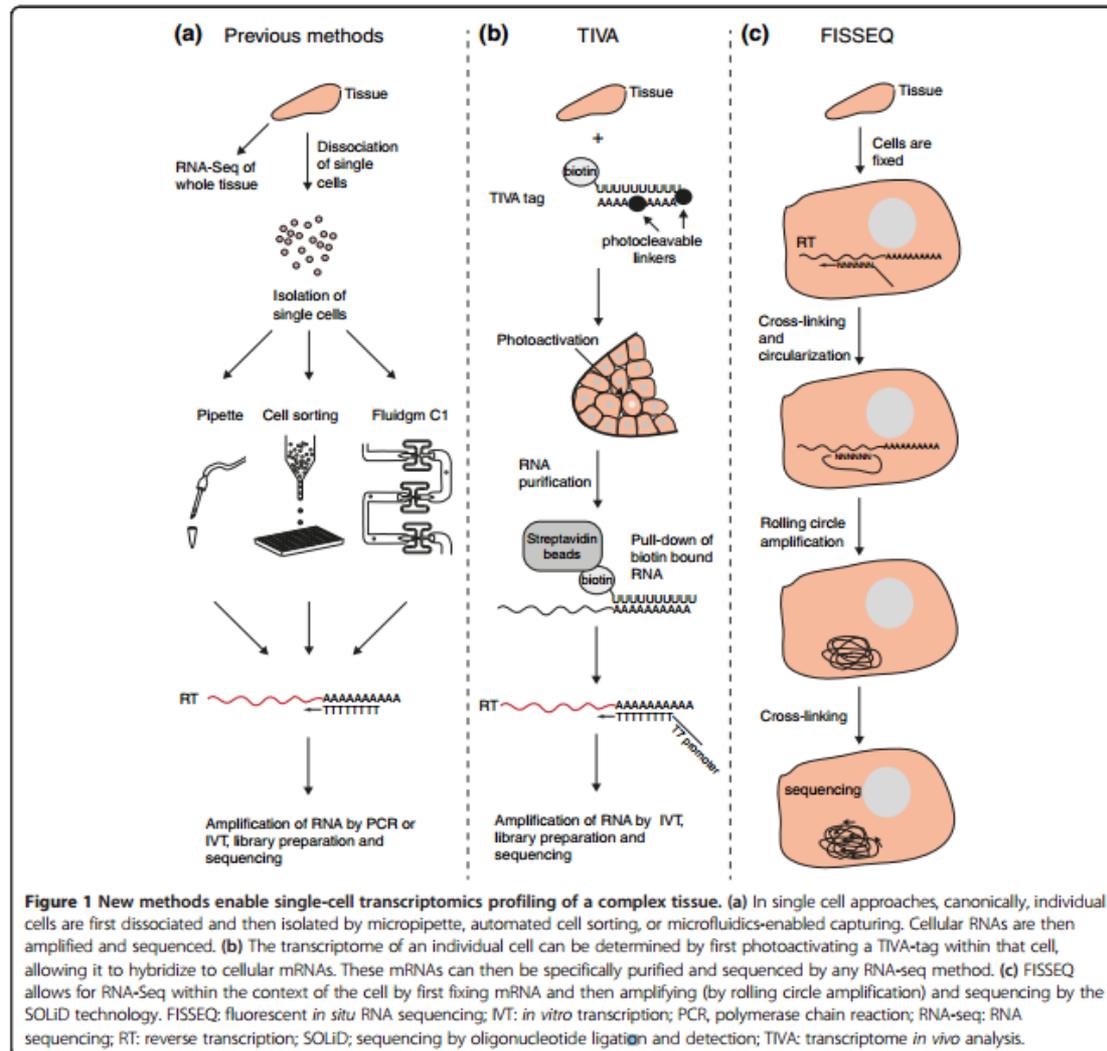
Alternative hypotheses on the origin of metastases.

# Seeing is believing: new methods for *in situ* single-cell transcriptomics

- ▶ Spatially defined single cells in live tissue
- ▶ Artifacts when separating by pipetting or using laser capture
- ▶ There is also a protocol for fixed cells in the tissue.



# *In situ* single-cell transcriptomics



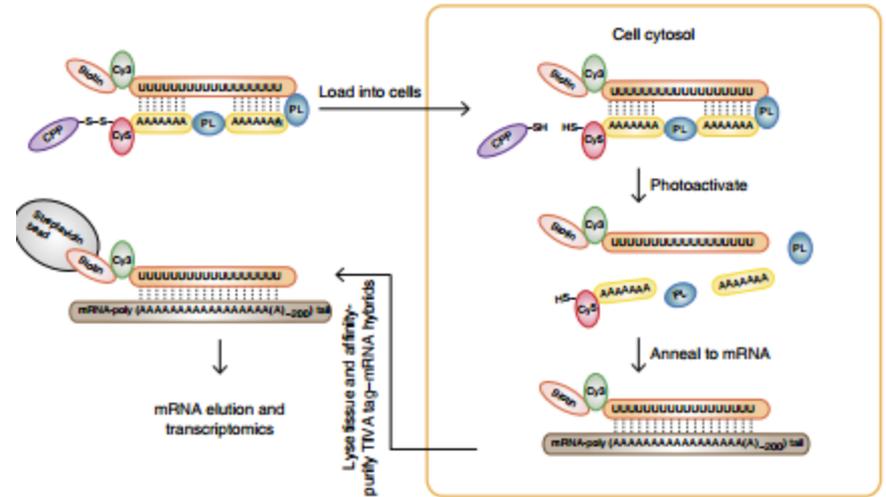
# Transcriptome *in vivo* analysis (TIVA)

TIVA tag:

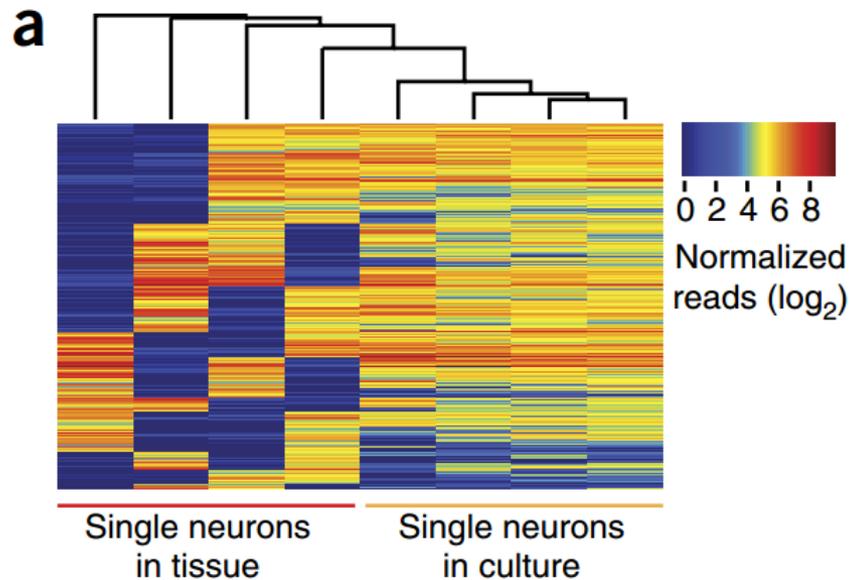
- a cell penetrating peptide,
- a photocleavable linker
- the fluorophores Cy3 and Cy5
- a poly(U) oligonucleotide
- biotin

TIVA tag:

- can permeate into the cells of a tissue the peptide dissociates from the TIVA tag
- laser photoactivation on a particular cell causes its TIVA tags to come undone
- tags can then anneal to mRNAs in the light-selected cell
- the desired mRNA can be pulled out by using streptavidin beads and then be sequenced



# Bimodal transcripts in single hippocampal neurons in tissue



(a) Heatmaps shows clustering of 645 bimodally expressed genes. Bimodally expressed genes were defined as having a gap in expression of at least four log units in two samples.

(b) Overlap between bimodal genes in single neurons from tissue and from culture (4 cells in each group).



**Thanks for your attention!**

**Any questions?**

