# Differential geometry and data science for single-cell biology
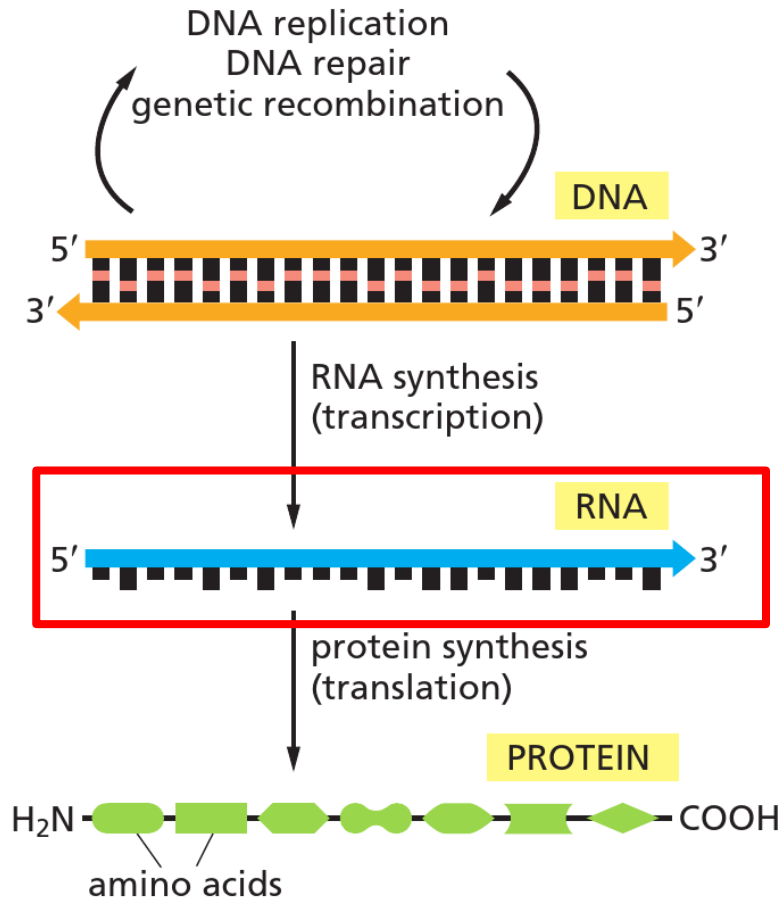
Ilia Kats and Luca Marconato

Stegle lab

dkfz.

GERMAN
CANCER RESEARCH CENTER
IN THE HELMHOLTZ ASSOCIATION

Research for a Life without Cancer

# Measuring gene expression


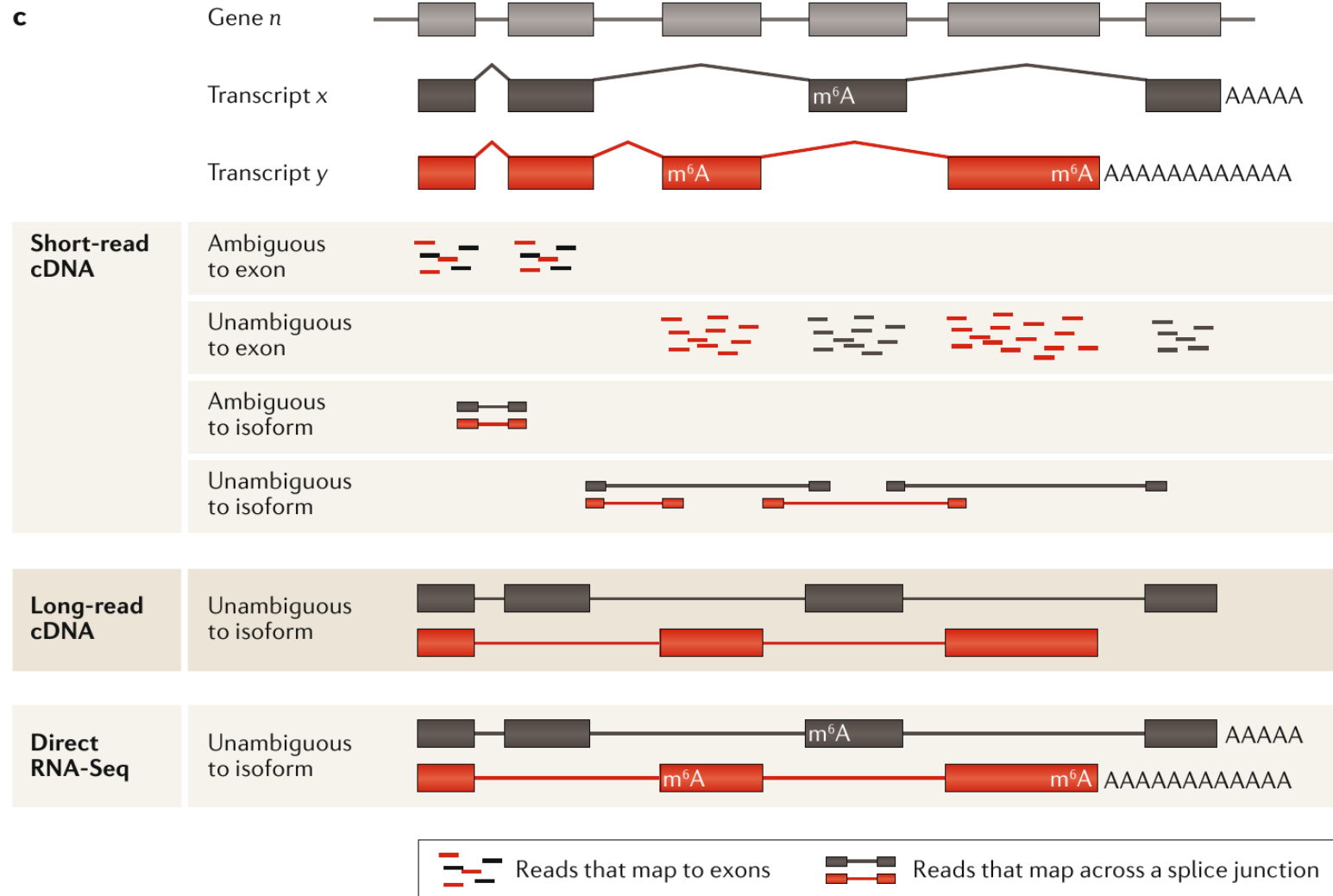
Gene expression: amount of gene product in the cell

Quantification by sequencing

Bruce Alberts, Molecular Biology of the Cell (6th edition, 2015)

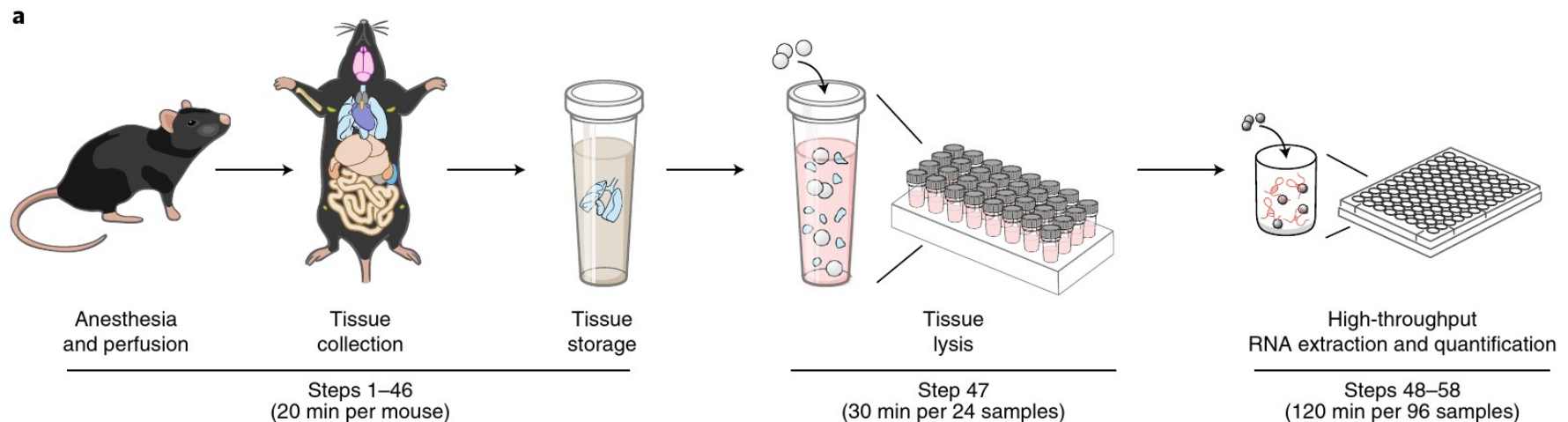dkfz.

# Measuring gene expression



Stark et al., 2019

# RNAseq: bulk vs. single-cell

- Bulk RNAseq
  - Works with homogenized tissues/tissue culture samples
  - Measures average expression in the sample
  - High sensitivity, low throughput
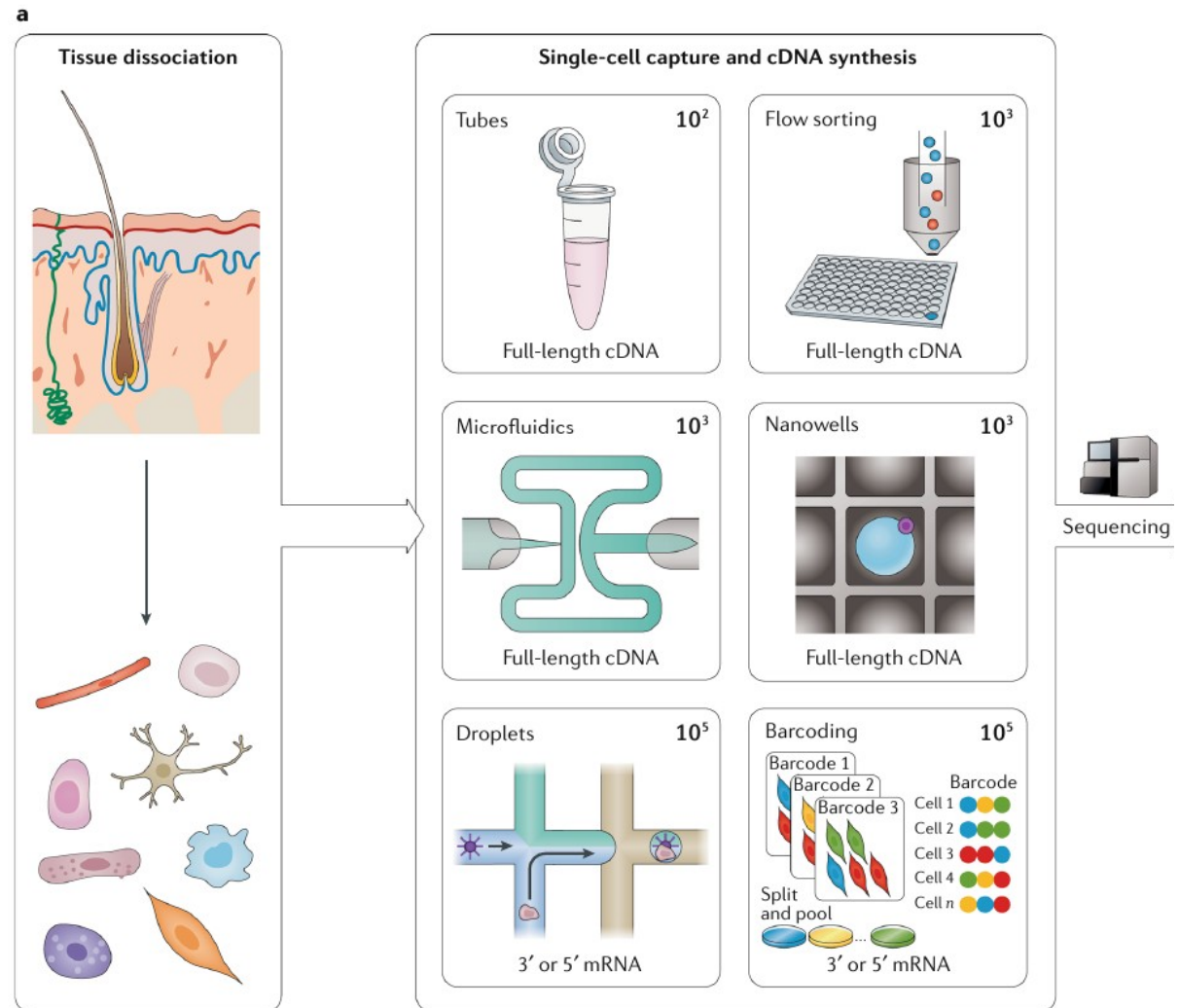  - Mostly useful for differential gene expression analysis



Pandey et al., 2020

# RNAseq: bulk vs. single-cell
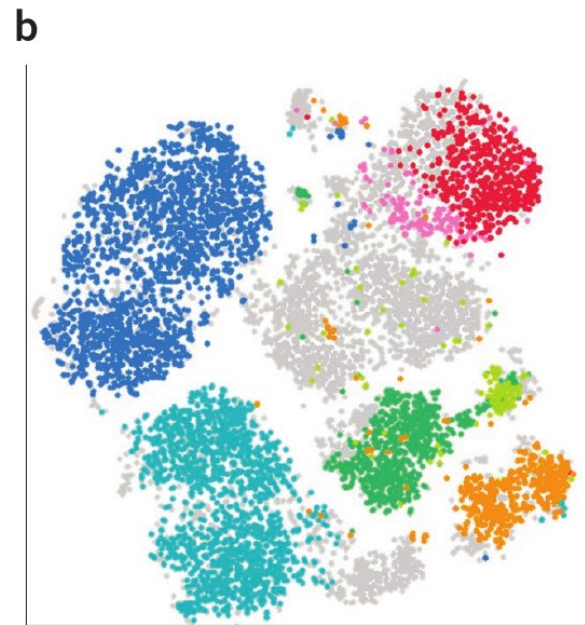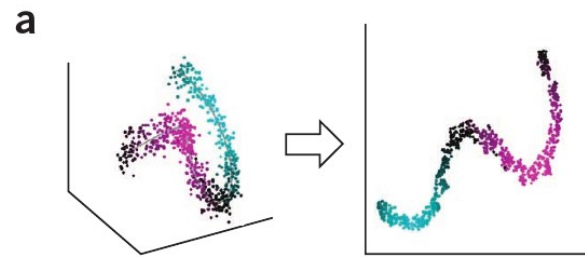
- **Single-cell RNAseq**
  - Measures gene expression in individual cells
  - Low sensitivity, high throughput
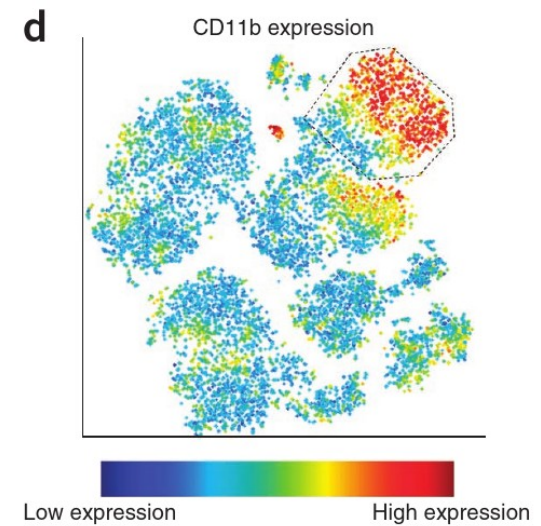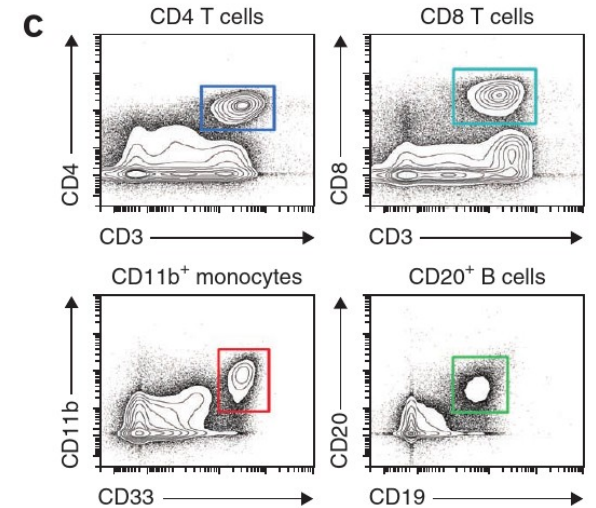  - Useful for: annotation of cell types, cell differentiation analysis, differential gene expression analysis, …



Stark et al., 2019

**dkfz.**

# Visualization: Dimensionality reduction

- One would like to plot the high-dimensional data in two dimensions

- Enables intuitive and easy visualization of clustering, batch correction, …

- Non-linear dimensionality reduction

  - tSNE

  - UMAP

  - PHATE

  - …



Amir et al., 2013

dkfz.

# tSNE (t-distributed stochastic neighbor embedding)

- Make the distribution of pairwise distances in 2 dimensions as similar as possible to the distance distribution in high-dimensional space → minimize KL-divergence

- Let $x_i$ be the coordinates of point i in high-dimensional space and $y_i$ the coordinates on low-dimensional visualization space

$$p_{i|j} := \frac{\exp\left(\frac{-\|x_i - x_j\|^2}{2\sigma^2}\right)}{\sum_{k \neq i} \exp\left(\frac{-\|x_k - x_i\|^2}{2\sigma^2}\right)}$$

Conditional probability in feature space (not symmetric)

$$\mathbb{P}_1 \text{ with } \mathbb{P}_1((i,j)) := \frac{p_{i|j} + p_{j|i}}{2n}$$

Symmetrized joint probability

$$\mathbb{P}_2 \text{ with } \mathbb{P}_2((i,j)) := q_{ij} := \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l}(1 + \|y_k - y_l\|^2)^{-1}}$$

Joint probability in visualization space

Hinton + van der Maaten, 2008

**dkfz.**

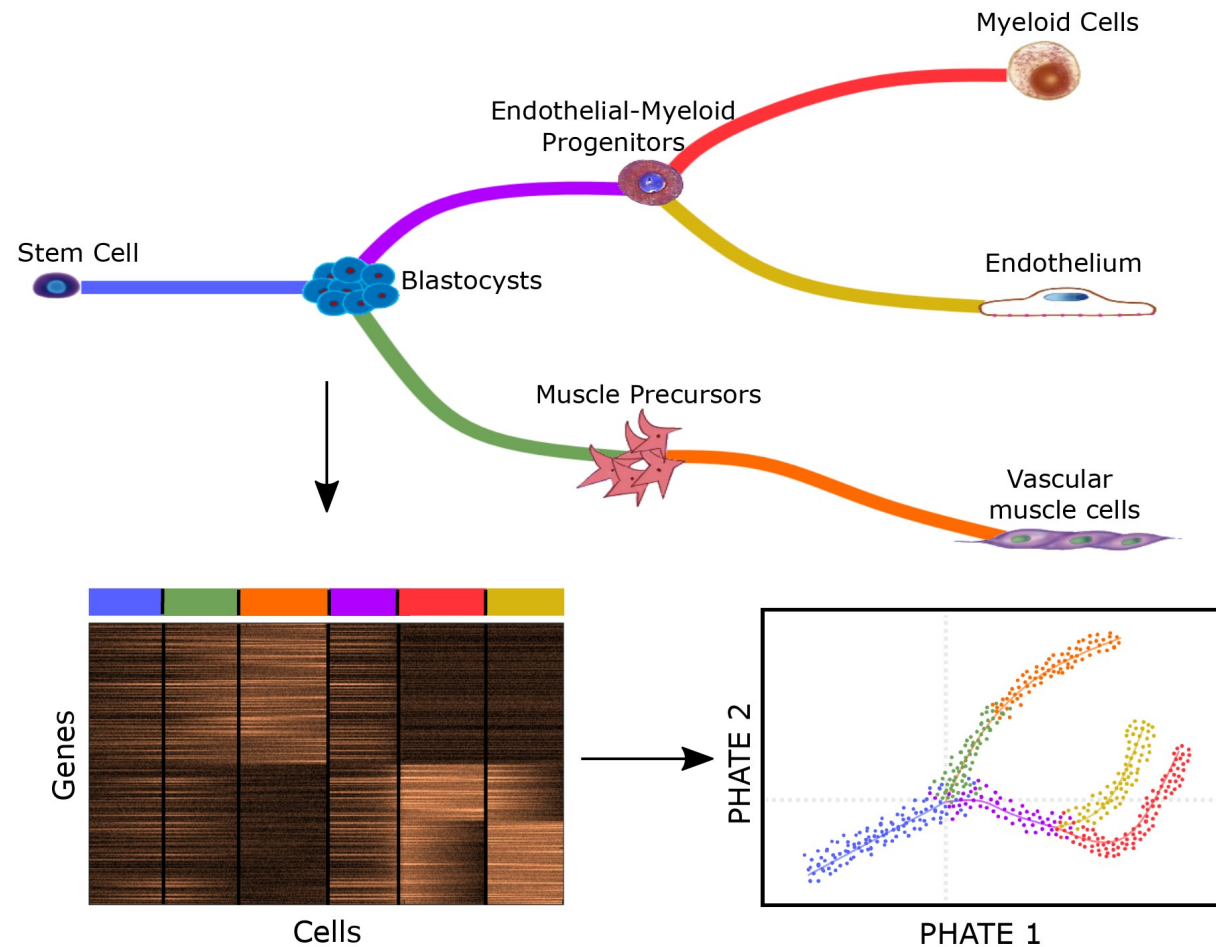# UMAP (Uniform Manifold Approximation and Projection)

- Approximate geodesic distances on data manifold by euclidean k-NN distances

- Use force-directed graph layout in visualization space

- Faster than tSNE, comparable embedding
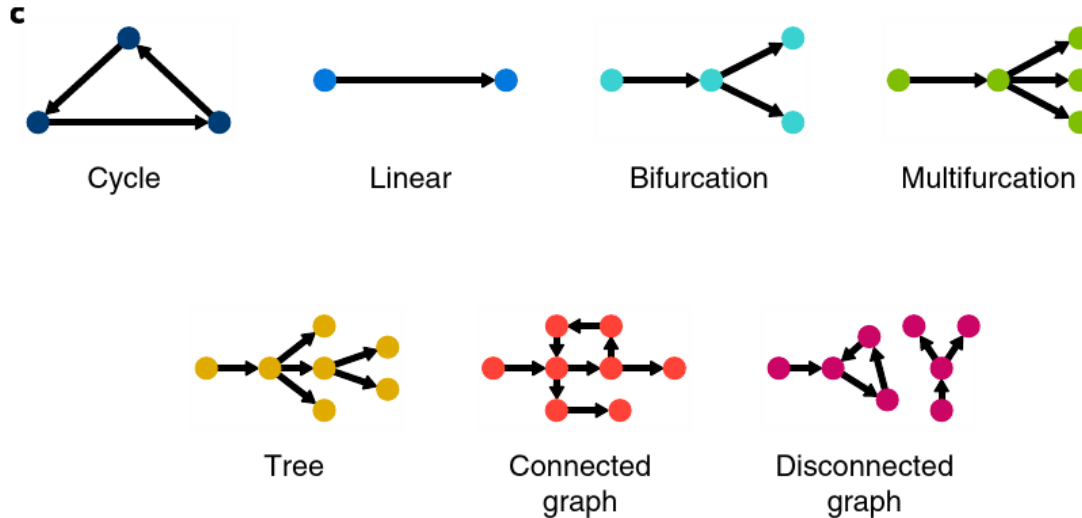
McInnes et al., 2018

**dkfz.**

# Trajectory inference

- Differentiating cells have a natural temporal ordering

- Differentiation typically not synchronized → cells from each stage are present in a sample

- Try to infer differentiation trajectory from a snapshot in time
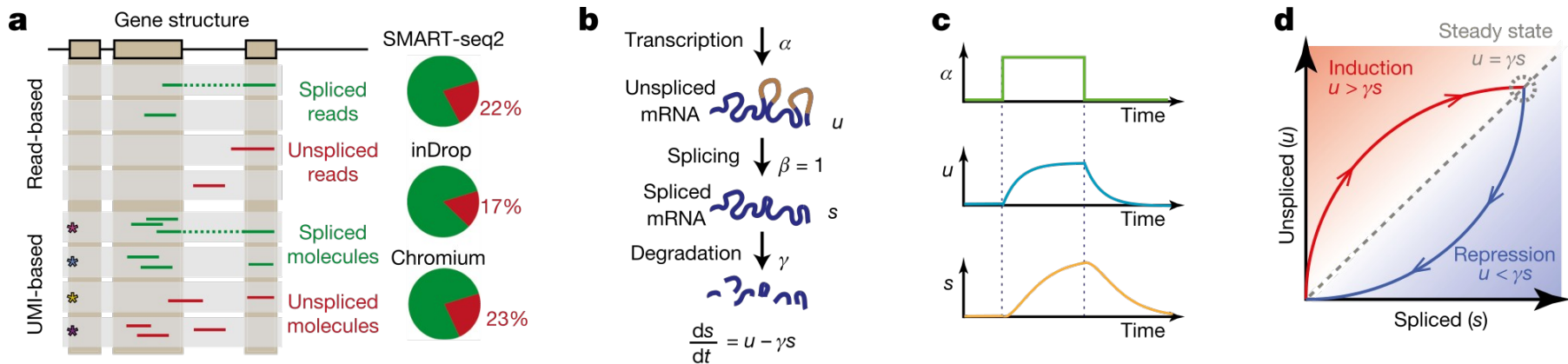


Moon et al., 2017

**dkfz.**

# Trajectory inference



Saelens et al., 2019

# RNA velocity

- It is possible to classify sequenced mRNAs into two classes
  - mature (old, being currently translated into protein)
  - Newly synthesized (will be translated in the future)
- This makes it possible to predict changes in gene expression for single cell (i.e. predict the future position of a cell on the manifold)
- Data extremely sparse



La Manno et al., 2018

**dkfz.**

# RNA velocity

- RNA velocity can be used for semi-automatic trajectory inference



Bergen et al., 2019

# Single-cell omics



Stuart + Satija, 2019

# MOFA+ (Multi-Omics Factor Analysis)

- Views: data modalities (e.g. RNAseq, protein abundance, …)

- Groups: batches/conditions



Argelaguet, Arnol, et al., 2019

**dkfz.**

# MOFA+ (Multi-Omics Factor Analysis)



Argelaguet, Arnol, et al., 2019

# Spatial omics



Stark et al., 2019

# Spatial omics



Eng et al., 2019

# Spatial omics: RNA velocity



Xia et al, 2019

# SpatialDE

- Models spatial gene expression as Gaussian Process

- Detection of spatially variable genes, classification into types of spatial variability



Svensson et al., 2018

# General issues in single cell omics

- Data sets are becoming larger very fast
  - Currently: largest available scRNAseq data set has approx. $10^6$ cells
  - Commercial platform for spatial RNAseq: 5000 spots
- Methods need to be computationally efficient and scale well
  - Increasing use of GPUs
- Very sparse data, large proportion of missing values
- No ground truth
  - Difficult to know if a method is doing The Right Thing™
- Data typically non-Gaussian
  - Makes exact calculations difficult, inefficient, or impossible
  - Sometimes transformation to approximate Gaussianity possible

**dkfz.**

# Microenvironments

- Core idea: interactions/communications between cells



Balkwill et al. The tumor microenvironment at a glance (J. Cell Sci 2012)

# Microenvironments



**a** IMC with 35 markers — pan-CK DNA, Fibronectin CK7, Ki-67 SMA — Cellular metaclusters — Neighbour interactions — Microenvironment communities

Adapted from Jackson et al. The single-cell pathology landscape of breast cancer (Nature 2020)



**a** Tissue sample — Number of cell–cell interactions — Matched randomized controls — Number of cell–cell interactions

**b** Row 7: is cell type X in the neighborhood of cell type 7? — Column 7: is cell type 7 in the neighborhood of cell type X?

**c** Cell phenotype of interest — Cell phenotype in neighborhood — Percent of significant images — Interaction / Avoidance

**d** Interaction / Avoidance — Normal — Grade 1 tumors — Grade 2 tumors — Grade 3 tumors — Clustered cell–cell interactions — Samples — Cluster 1 — Cluster 2

**e** Cluster 1 — Cluster 2 — Interaction / Avoidance

Schapiro et al. histoCAT: analysis of cell phenotypes and interactions in multiplex image cytometry data (Nat. methods 2017)

## Approach:

1) Instance segmentation

2) Graph construction, features averaging pixel signal

3) Permutation test

4) Clustering

dkfz.

# Microenvironments



Applying the method presented in Carpenter et al. (Nat. methods, 2018)



**Multi-gigapixel WSI**

**Patches**

**Slide Graph**

**Conventional**

WSI Level Graph Representation

Convolutional Layer

Pool

Dense

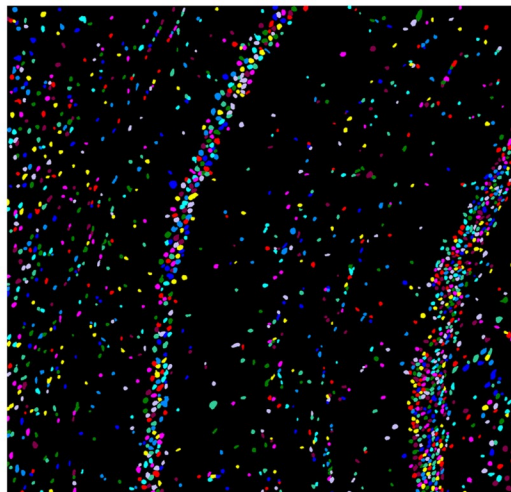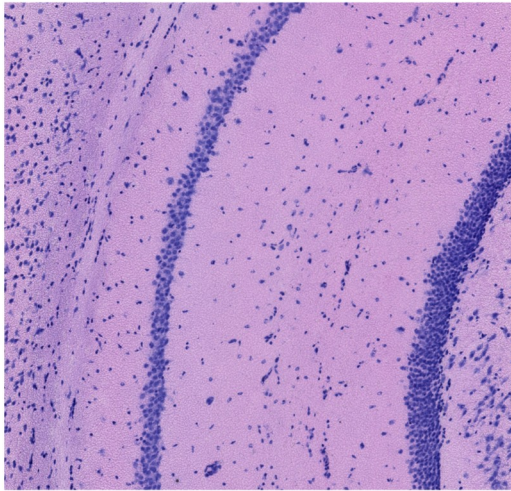Patch Level Image Classifier

GCN

Patch Score Aggregation

**WSI-level Prediction**

**WSI-level Prediction**

## WSI = Whole slide image

Lu et al. Capturing Cellular Topology in Multi-Gigapixel Pathology Images (CVPR 2020)

Approach:

1) Instance segmentation from H&E stained images

2) Graph construction, features by clustering image features

3) Predictions with GCN

# Alignment

- Same biology, different experiments → need to match the data



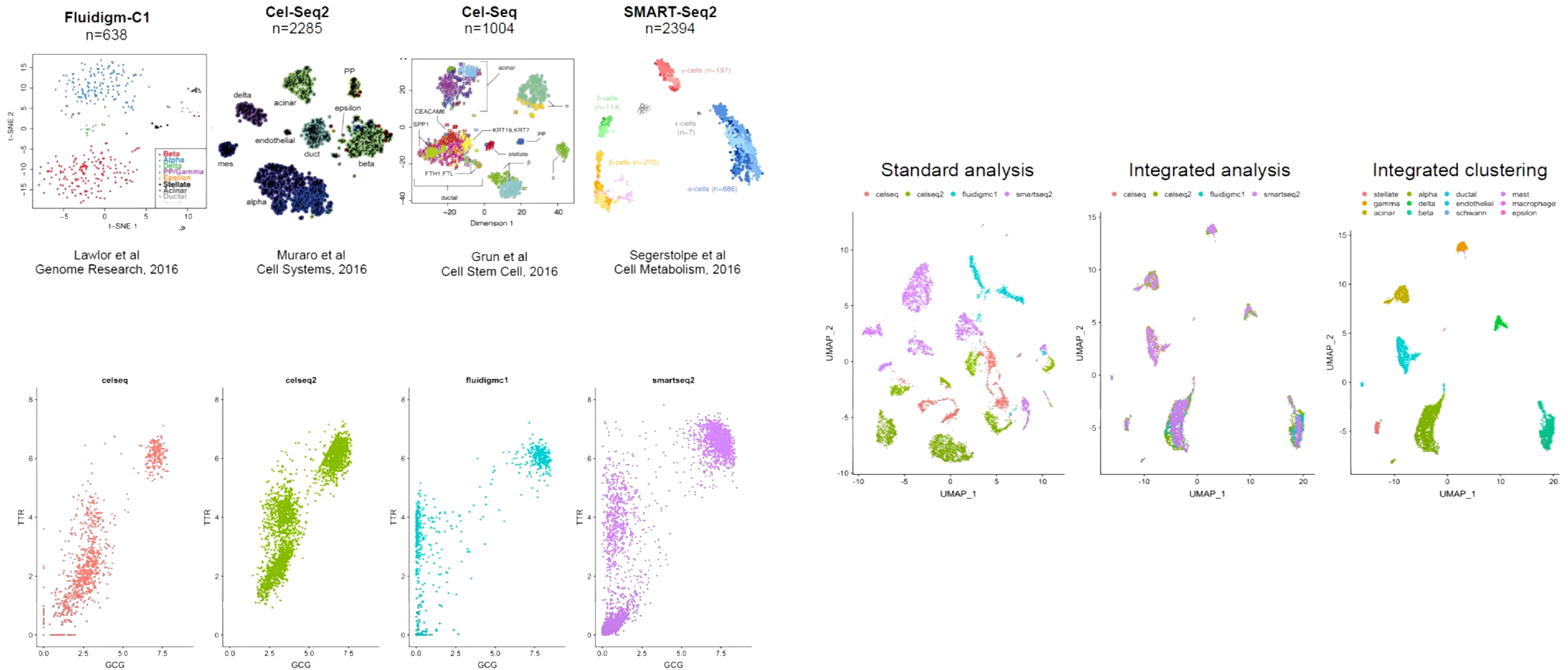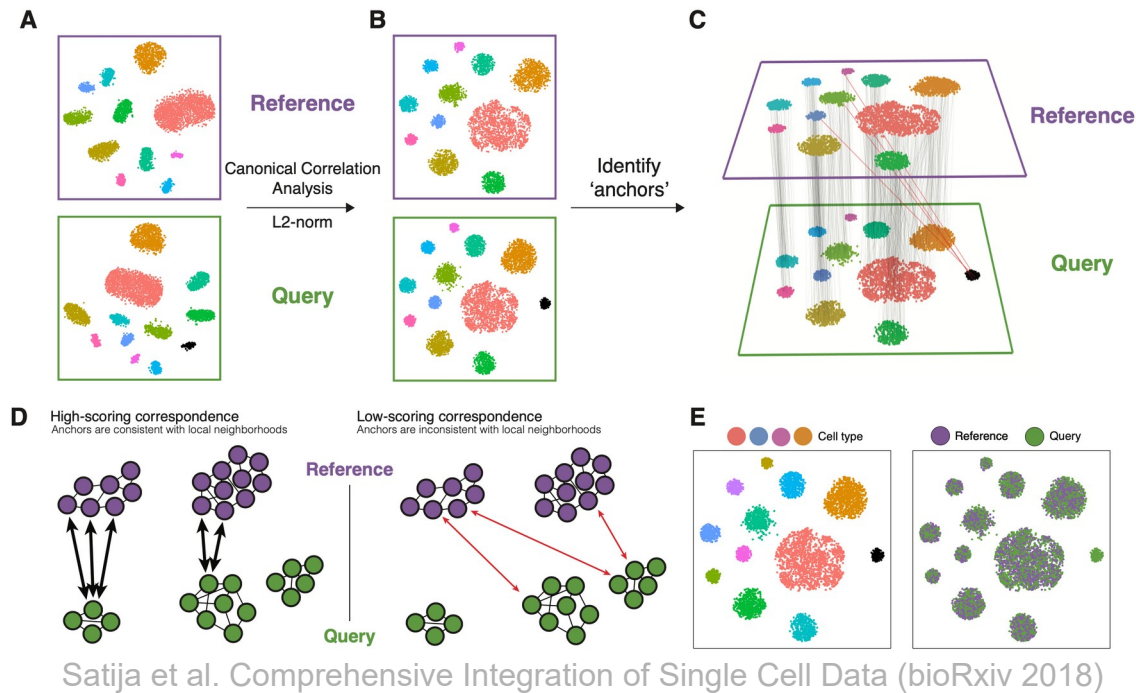Adapted from Comprehensive Integration of Single Cell Data (presentation of Rahul Satija)

# Alignment



Satija et al. Comprehensive Integration of Single Cell Data (bioRxiv 2018)

Approach:

1) Joint dimensionality reduction

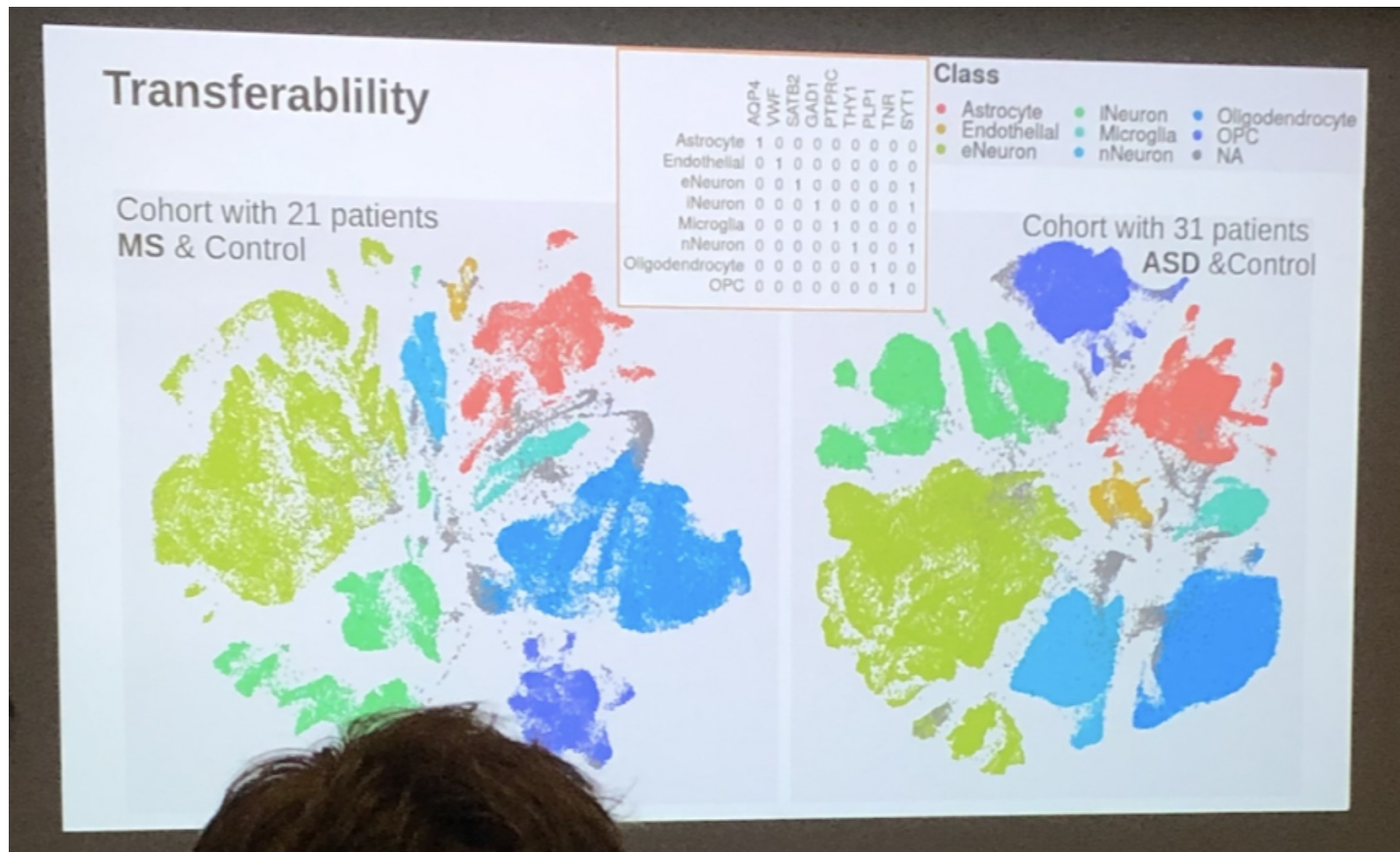$$X_{f,c}, Y_{f,c} \quad f \text{ features}, c \text{ cells, same features}$$

$$\max_{u,v} u^\top X^\top Y v, \quad \|u\|^2 \leq 1, \|v\|^2 \leq 1$$

2) Finding anchors by mutual k nearest neighbors

3) Anchor weighting (by a Gaussian kernel) and alignment

dkfz.

# Alignment

- Idea: considering the topology for improving the alignment



Picture from a presentation of F. Frauhammer, DKFZ

# Some directions that could be worth exploring

- Topology-aware manifold aligner

- Differential expression in manifolds

- Geometric deep learning

Datasets:

- Many modalities
  - Paired (many modalities for the same entity)
  - Unpaired (from the same or from different samples)

- High dimensionality (>20000 genes in humans)

- Multi-channel, high-resolution images

- Data with tree-like latent structure (suitable for hyperbolic embeddings)

**dkfz.**

**Thank you
for your attention!**

**dkfz.** GERMAN
CANCER RESEARCH CENTER
IN THE HELMHOLTZ ASSOCIATION

Research for a Life without Cancer

# SpatialDE

- Boundary effects?

- Euclidean distance smaller than geodesic distance

- Difference in detected patterns?

**dkfz.**