

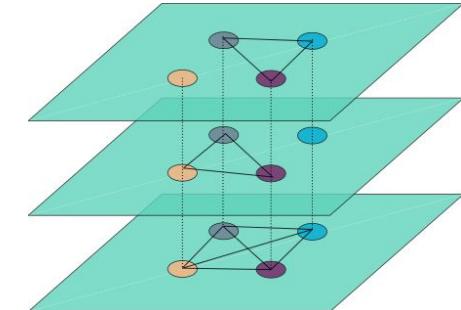
~~Clustering~~ ^{UNSUPERVISED LEARNING}

in Computational Biology

Karolina Sienkiewicz

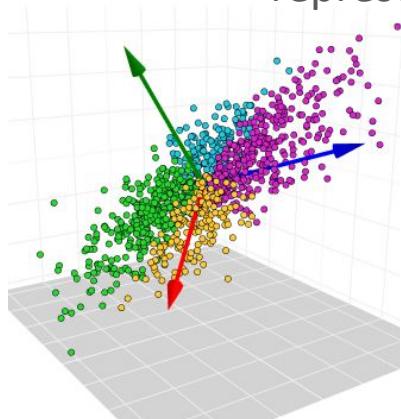
Weill Cornell Medicine, Cornell University

sienkiewicz2k@gmail.com  @sienkieee



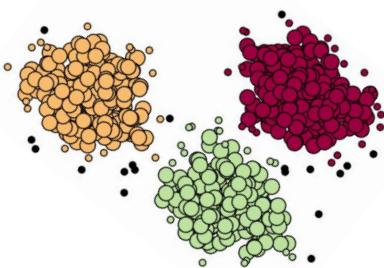
Unsupervised learning

low-rank
representation



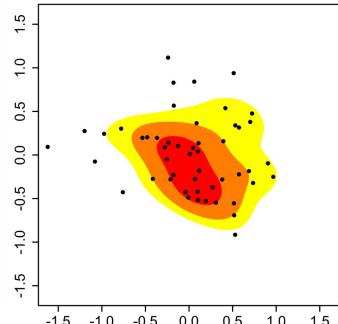
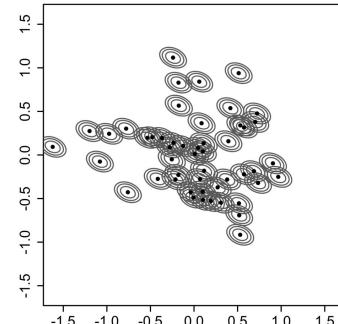
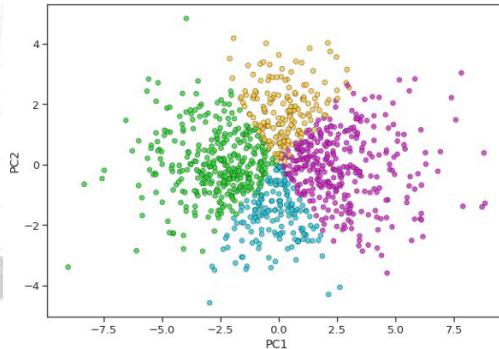
From: Casey Chang

clusters



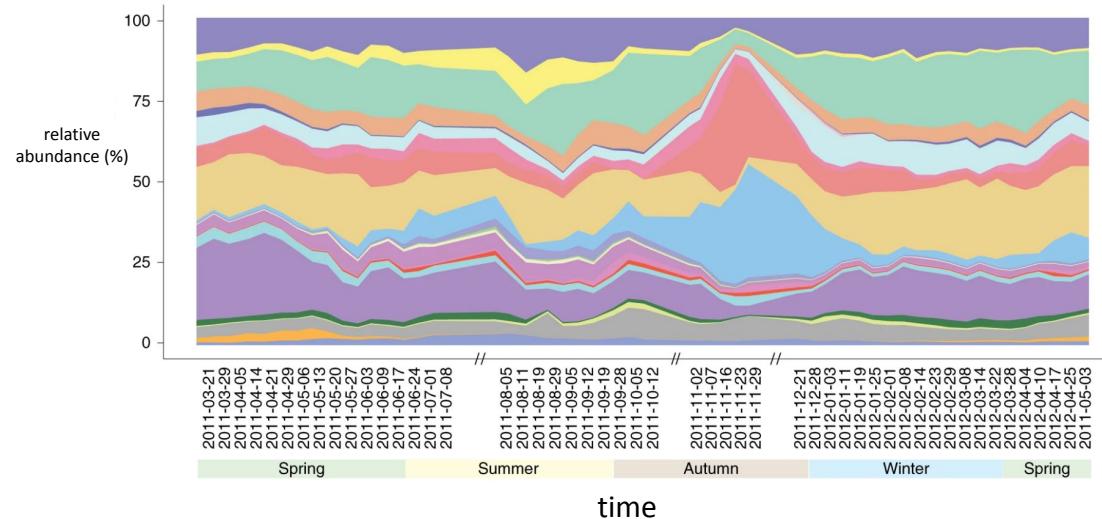
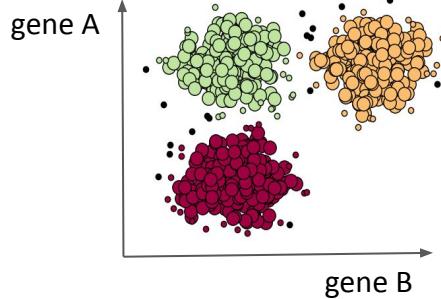
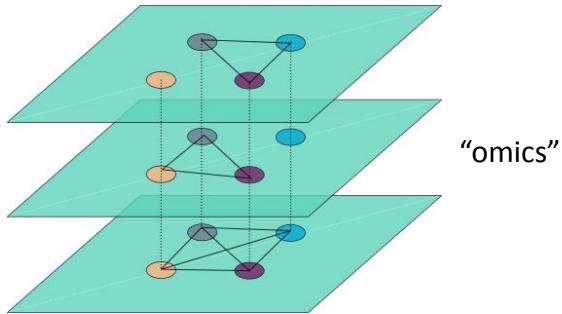
probability
density
estimation

$$P \approx P_{\text{data}} \text{ model}$$

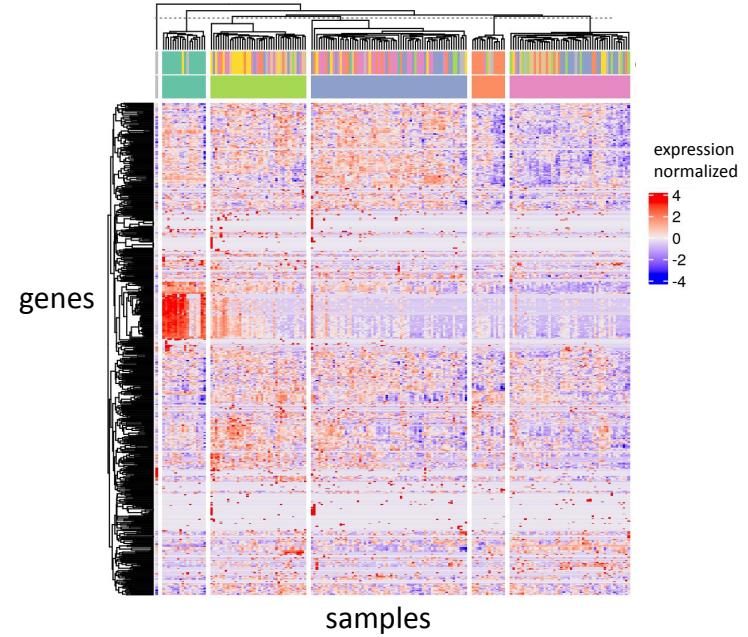


wikiwand.com/en/Multivariate_kernel_density_estimation

Biological datasets are complex



PMID:
33139880



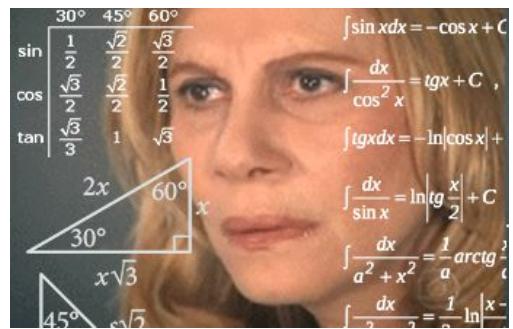
expression normalized
4
2
0
-2
-4

UNSUPERVISED LEARNING

Practical application #1:
Batch effect correction

Batch effect = non-biological variation

Why is there
a batch effect
in my data?

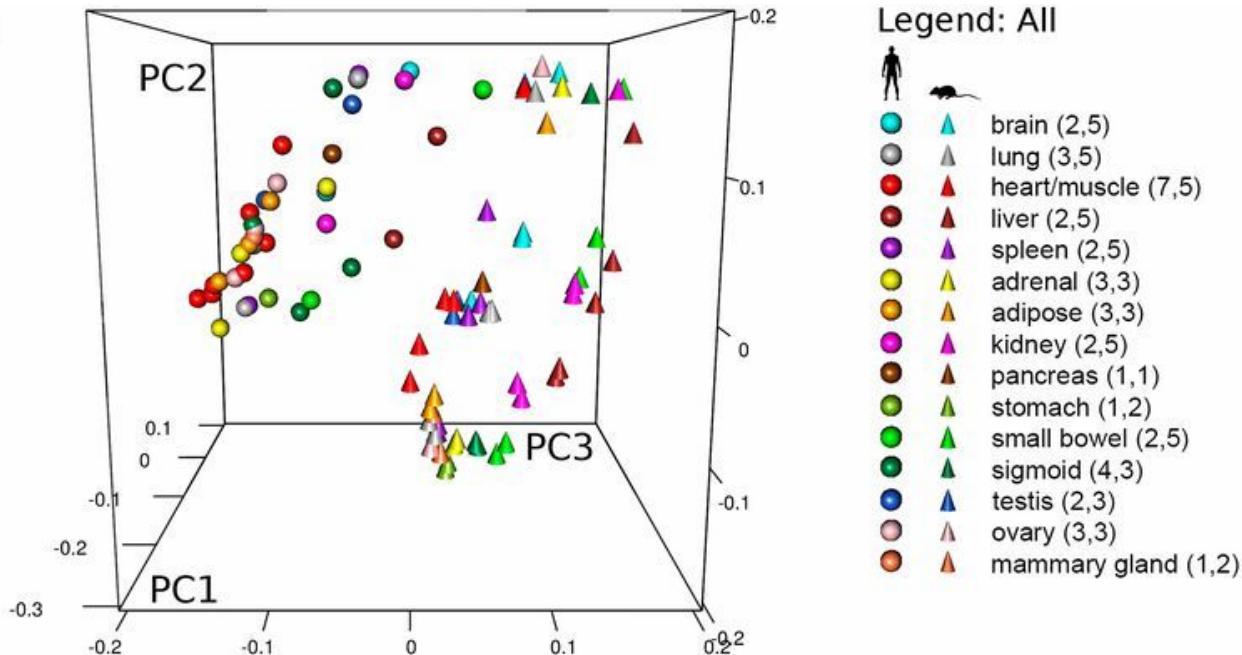


- Different days/months of the experiments
- Different reagents (enzymes, buffers)
- Different mice (from different companies)
- Different sequencers
- Lab protocol or experimenter/technicians

Batch effect example

“Gene Expression Is More Similar Among Tissues Within a Species Than Between Corresponding Tissues of the Two Species”

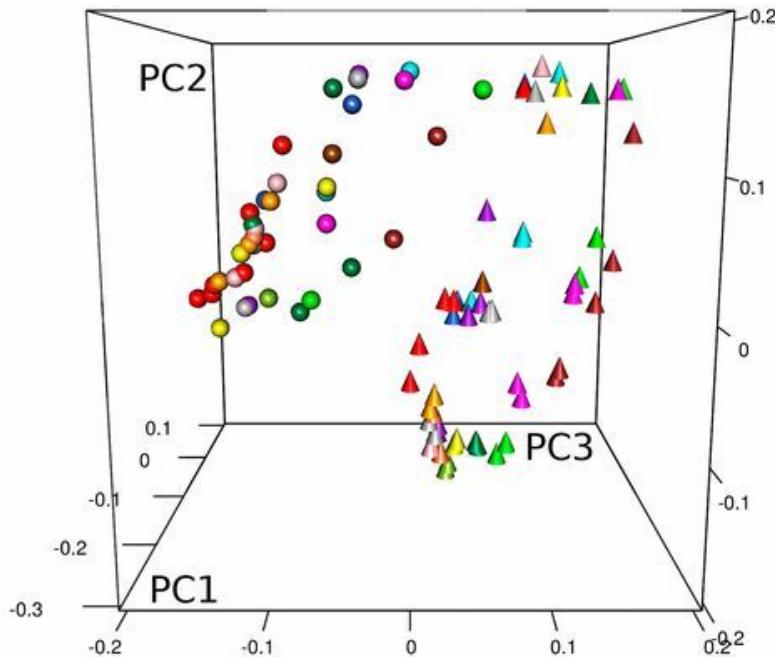
A



Batch effect example

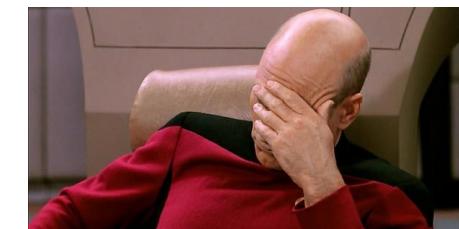
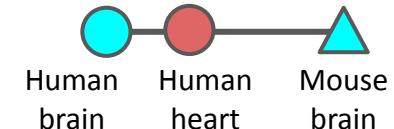
“Gene Expression Is More Similar Among Tissues Within a Species Than Between Corresponding Tissues of the Two Species”

A



Legend: All

- brain (2,5)
- lung (3,5)
- heart/muscle (7,5)
- liver (2,5)
- spleen (2,5)
- adrenal (3,3)
- adipose (3,3)
- kidney (2,5)
- pancreas (1,1)
- stomach (1,2)
- small bowel (2,5)
- sigmoid (4,3)
- testis (2,3)
- ovary (3,3)
- mammary gland (1,2)

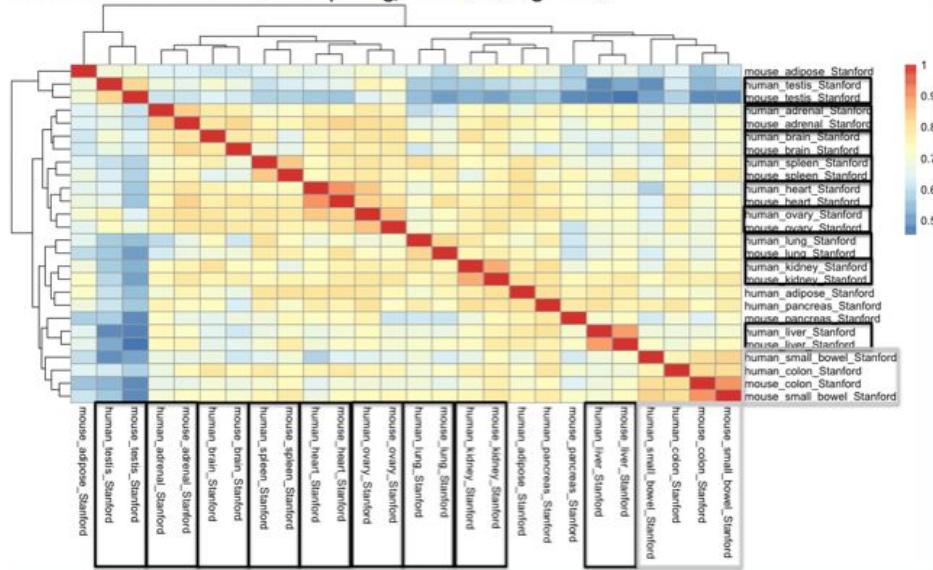
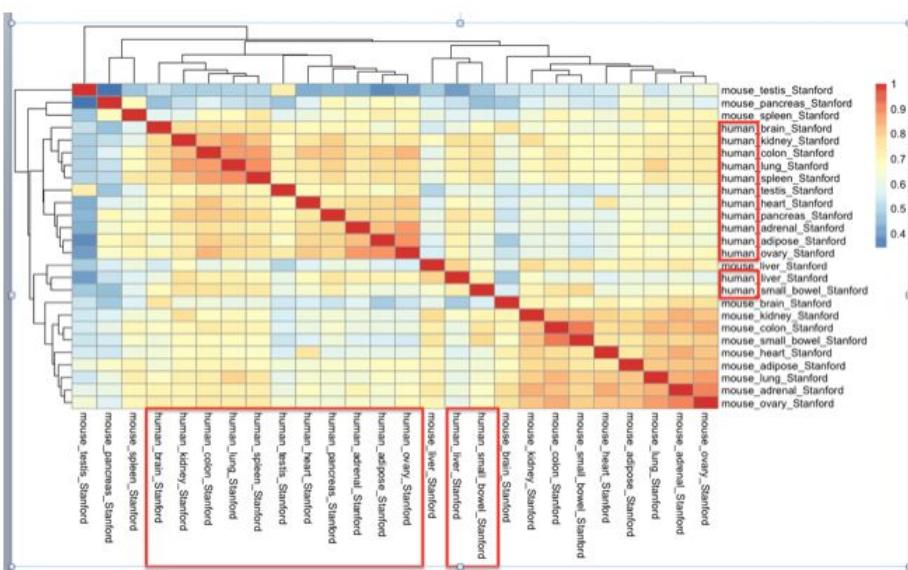




Yoav Gilad
@Y_Gilad

We reanalyzed the data from pnas.org/content/111/48... and found the following:

The original analysis in the paper, considering only the samples that were sequenced at Stanford (data cluster by species):



by Yoav Gilad & Orna Man

Good experimental design minimizes batch effects

consistent
pre-processing
of all samples

keeping track
of metadata

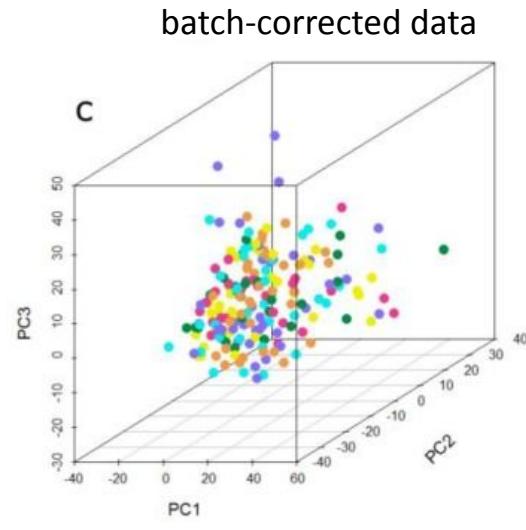
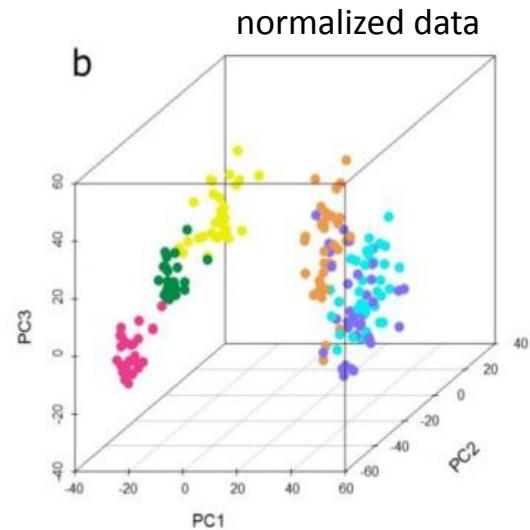
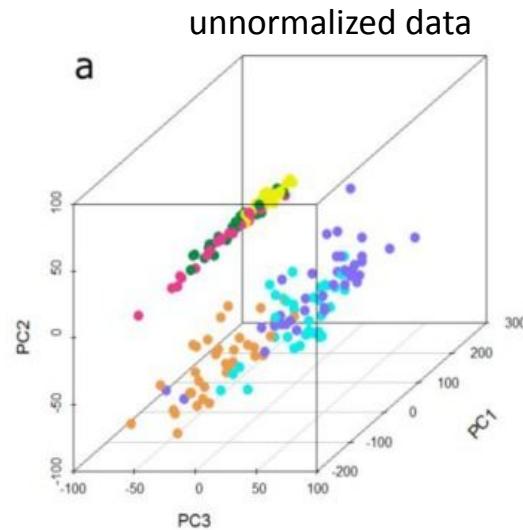
balancing
groups of interest

avoiding
“perfect confounding”



Combating batch effect

visualize/cluster → identify → correct



ComBat

Johnson et al (2007) *Biostatistics*

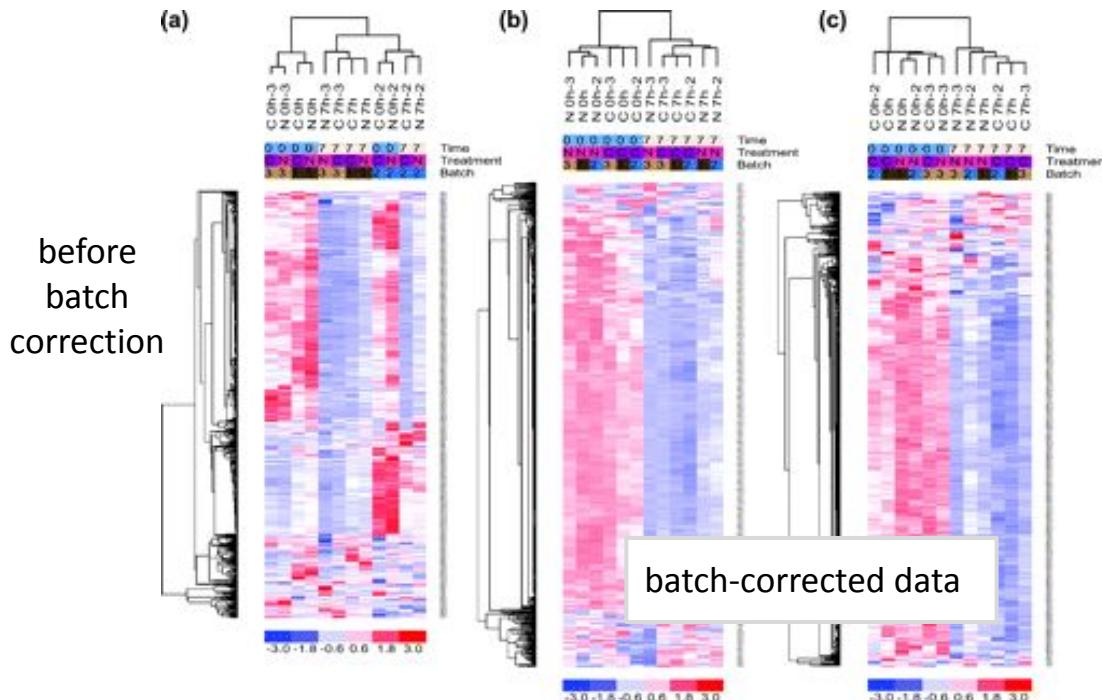
limma

complex (e.g. nested) batches

Brezina et al (2015) *Microarray*

Combating batch effect

visualize/cluster → identify → correct



ComBat

Johnson et al (2007) *Biostatistics*

limma

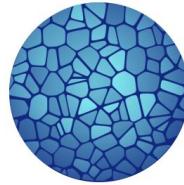
complex (e.g. nested) batches

Johnson et al (2007) *Biostatistics*

UNSUPERVISED LEARNING

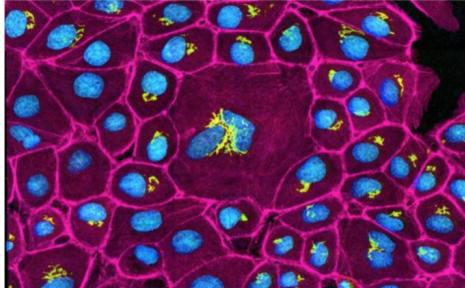
Practical application #2-4:
Single-cell gene expression analysis

Single-cell technologies - motivation

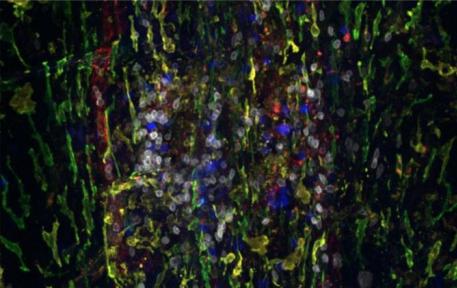


HUMAN
CELL
ATLAS

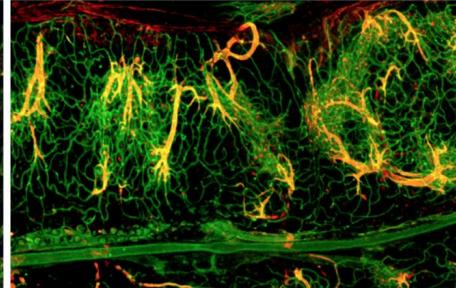
Skin epithelium



Brain meninges



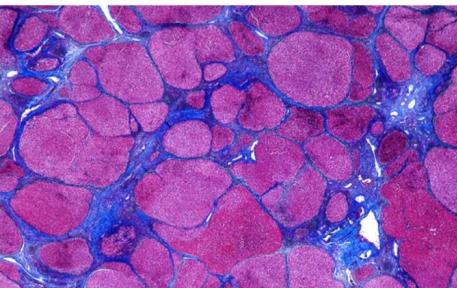
Blood vessels



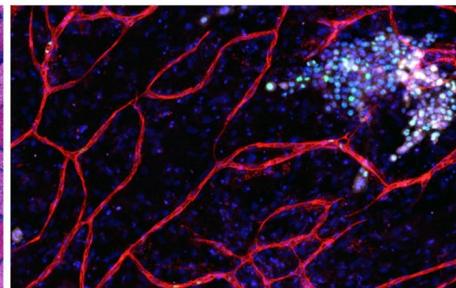
Small intestine



Liver cirrhosis



Breast cancer



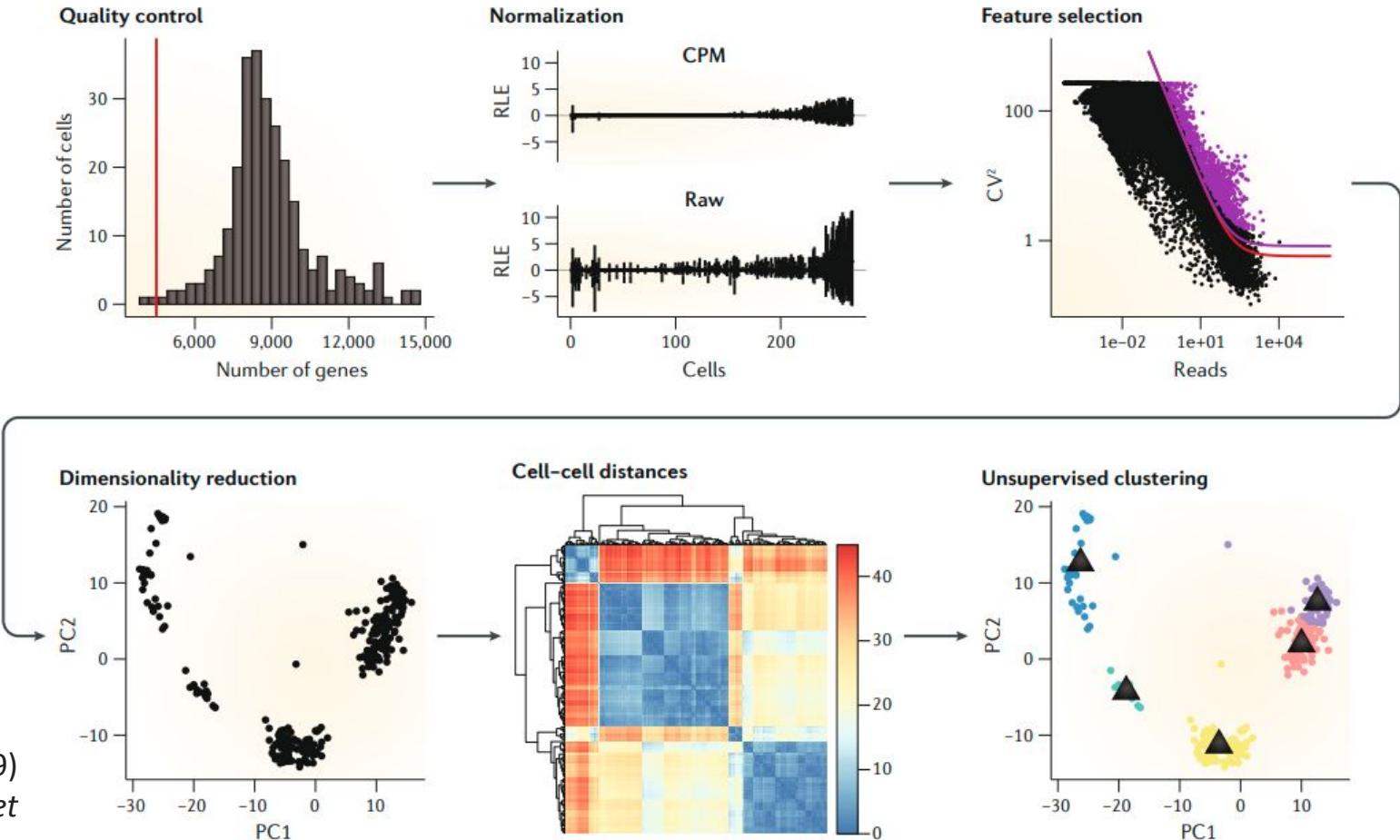
10x Genomics
(2016)

>1.3 million
cells

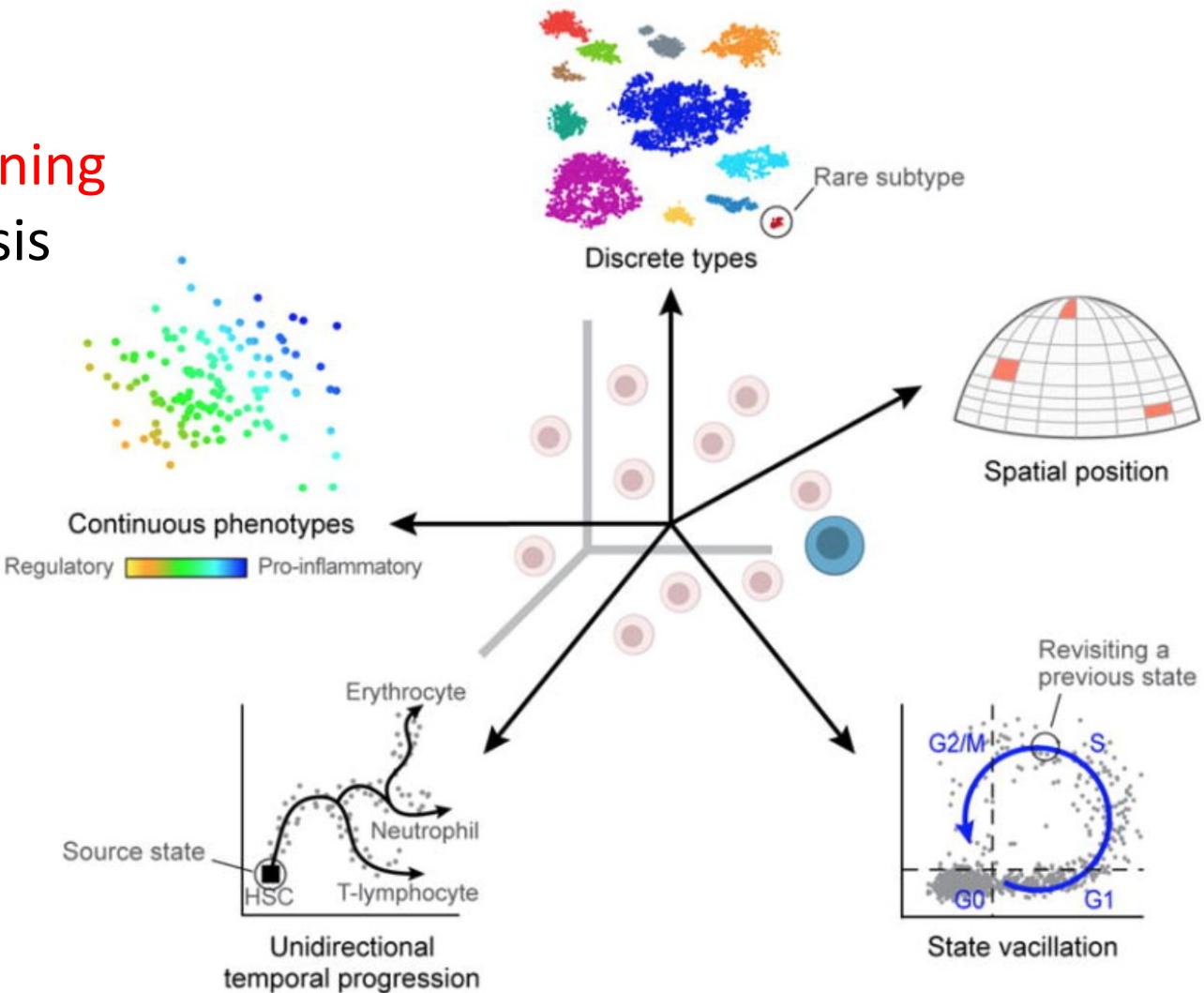
Tang et al. (2009)
Nat Methods

8 cells

Single-cell analysis common workflow



Many faces of Unsupervised Learning in single-cell analysis

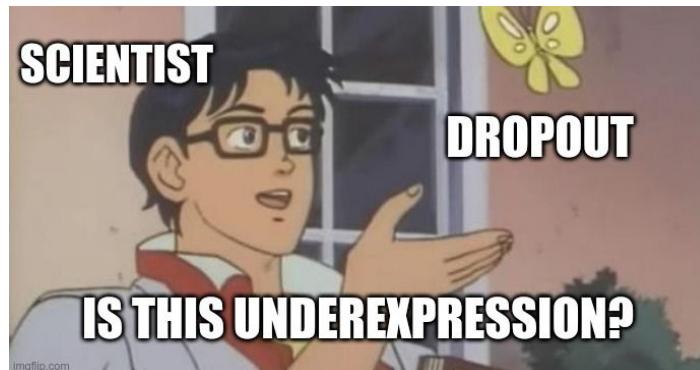


Single-cell technologies - **technical** considerations

cell-specific capture efficiency

amplification bias

library quality



Dropout in single-cell

stochastic transcription

low abundance of transcripts

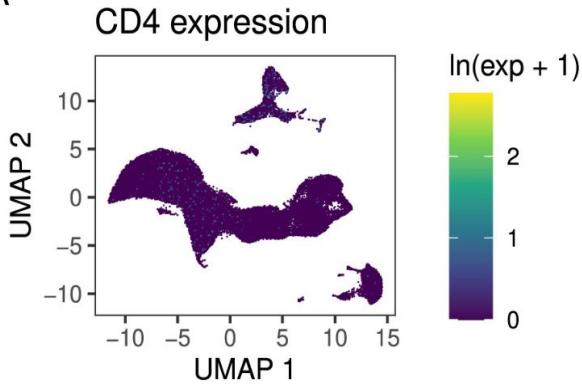
shallow sequence depth

Kernel density estimation example

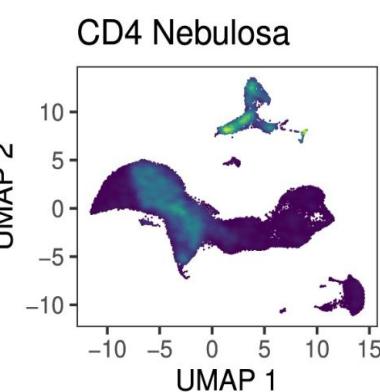
Nebulosa
inference based on nearby cells

$$\hat{f}(x; H) = \frac{1}{n} \sum_{i=1}^n w_i K_H(x - X_i)$$

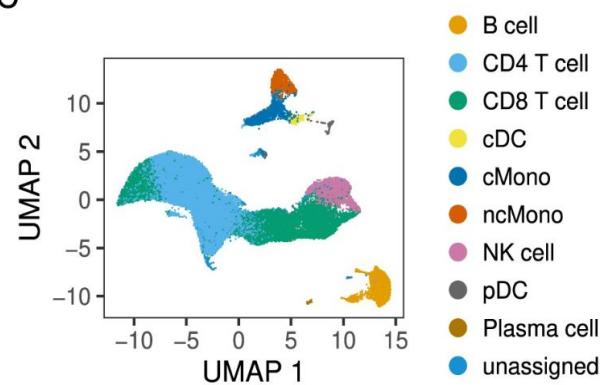
A



B



C

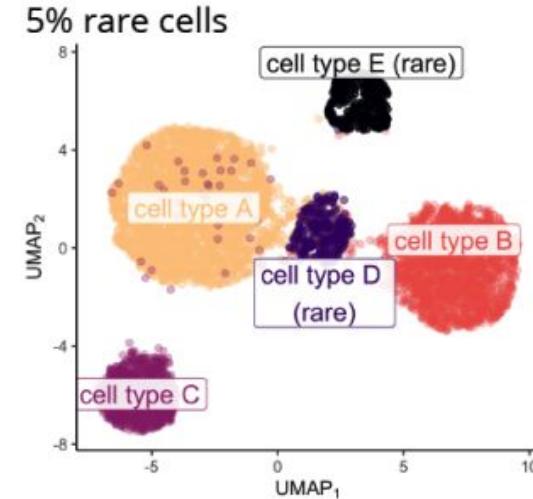
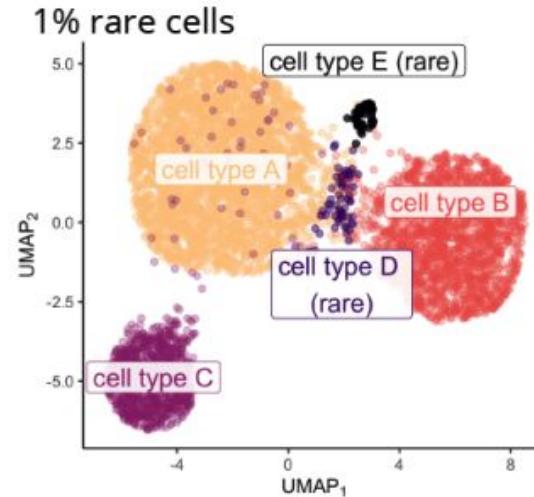
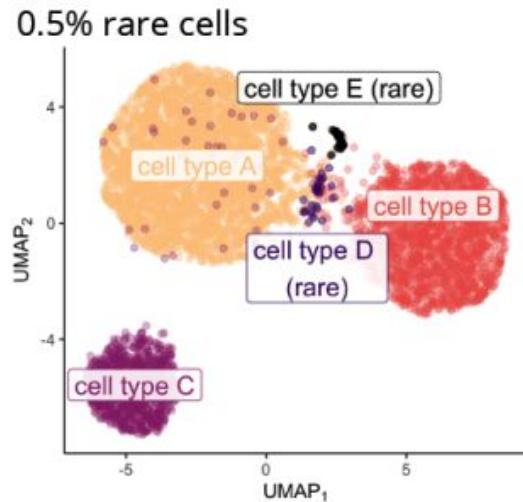


Identifying rare clusters

Detecting rare clusters is challenging

Increasingly rare cell types

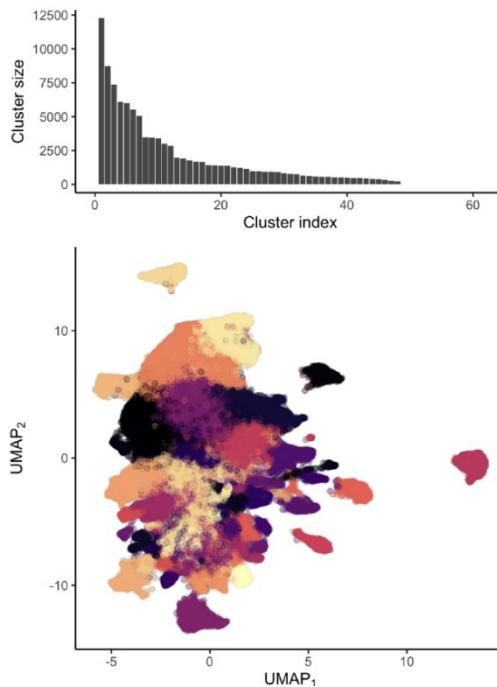
(A) Ground truth clusters (varying % of rare cells)



Rarity: Discovering rare cell populations from single-cell imaging data
<https://www.biorxiv.org/content/10.1101/2022.07.15.500256v1>

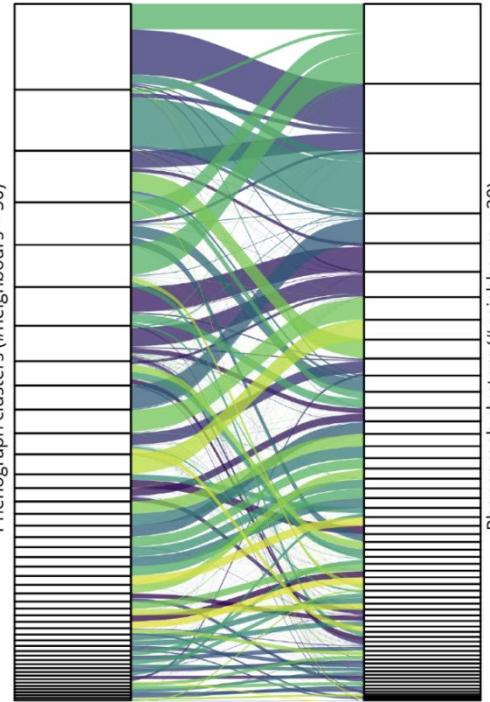
Can we simply increase k to increase granularity of clusters?

A Phenograph clusters (#neighbours = 50)

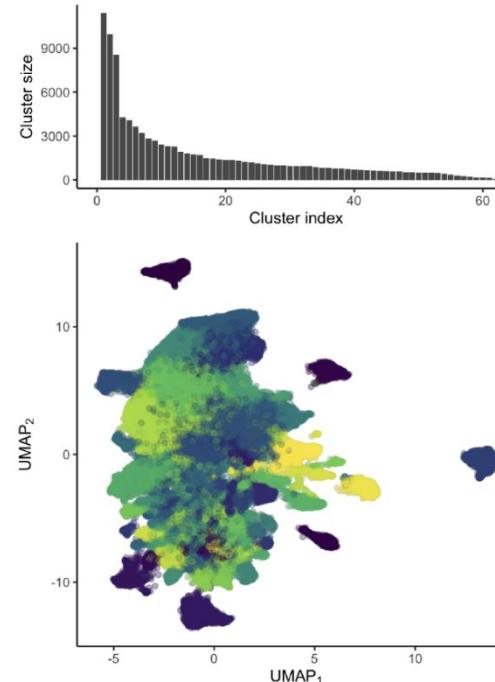


Increasing the number of clusters

Phenograph clusters (#neighbours = 50)

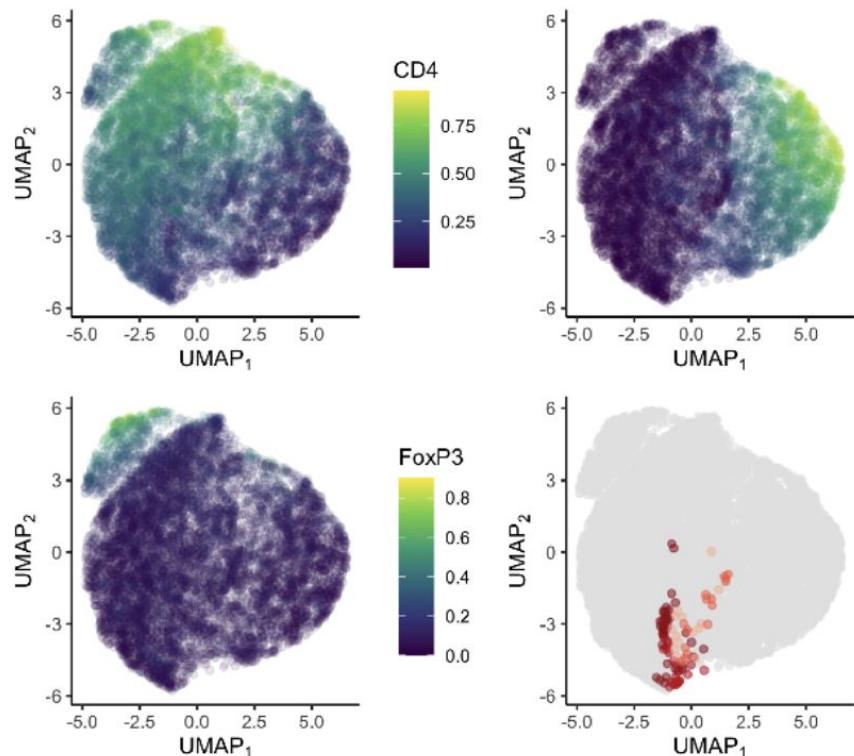


Phenograph clusters (#neighbours = 20)

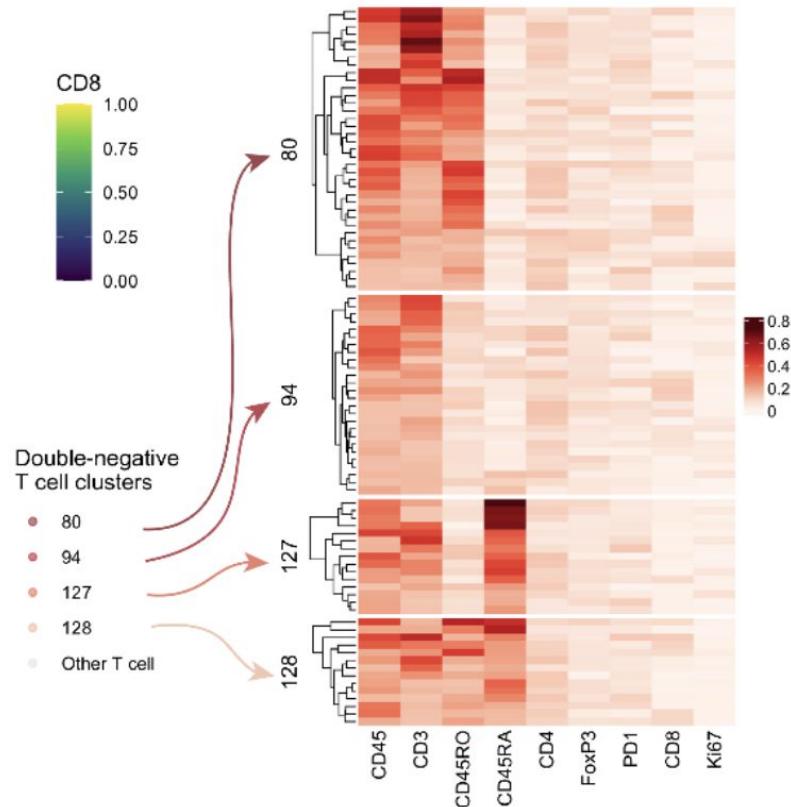


Example: discovering rare CD4- CD8- T cell clusters

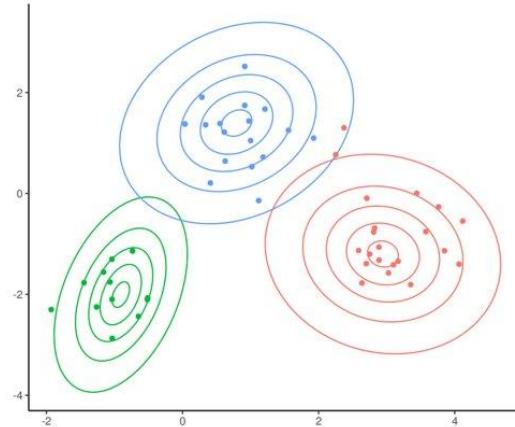
C Focusing on T cells



D Double-negative T cell clusters



The everlasting question: How to choose “k” in clustering models



Can we circumvent choosing the number of clusters?

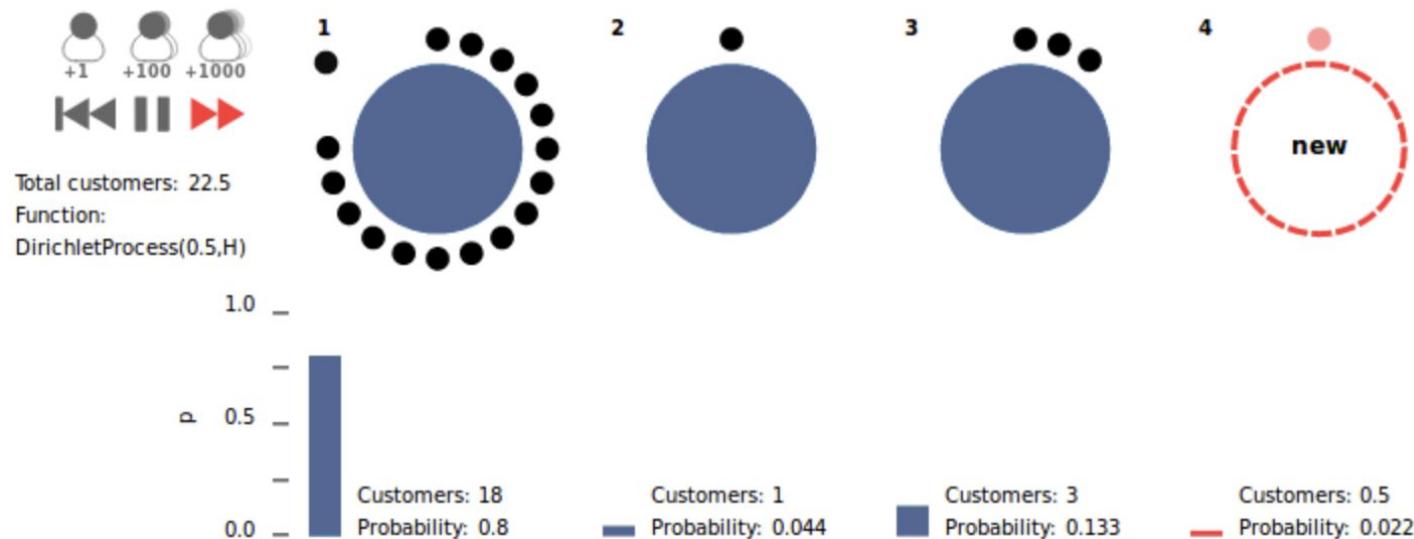
To some extent: there is an entire research field “**Bayesian nonparametrics**” trying to circumvent the issue by allowing a flexible (potentially infinite) number of clusters

Desirable properties of Bayesian nonparametric models:

- When new data comes in, ***new clusters*** should be ***allowed to appear*** (i.e. the number of clusters should not be fixed)
- We might want to capture ***uncertainty*** about the number of clusters

Flexible number of clusters with the Chinese Restaurant Process

We could use a Gaussian Mixture Model, where the number of clusters is specified by the Chinese Restaurant Process (thus allowed to be flexible)



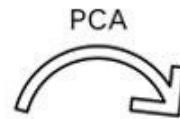
See gif in [https://commons.wikimedia.org/wiki/File:Chinese_Restaurant_Process_for_DP\(0.5,H\).gif](https://commons.wikimedia.org/wiki/File:Chinese_Restaurant_Process_for_DP(0.5,H).gif)

Probabilistic and non-linear extensions of PCA

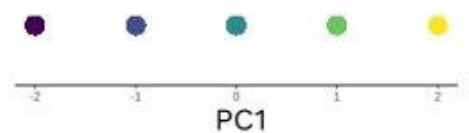
Probabilistic PCA: turning PCA into a generative model

PCA

Observed data



Low-dimensional latent representation

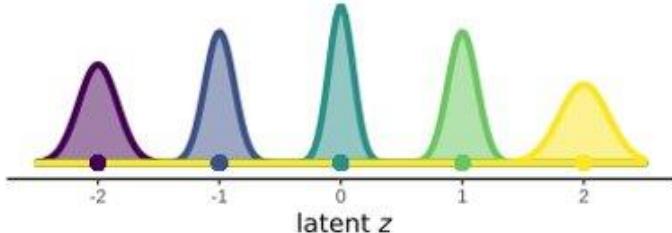


Probabilistic PCA

Observed data



Low-dimensional latent representation

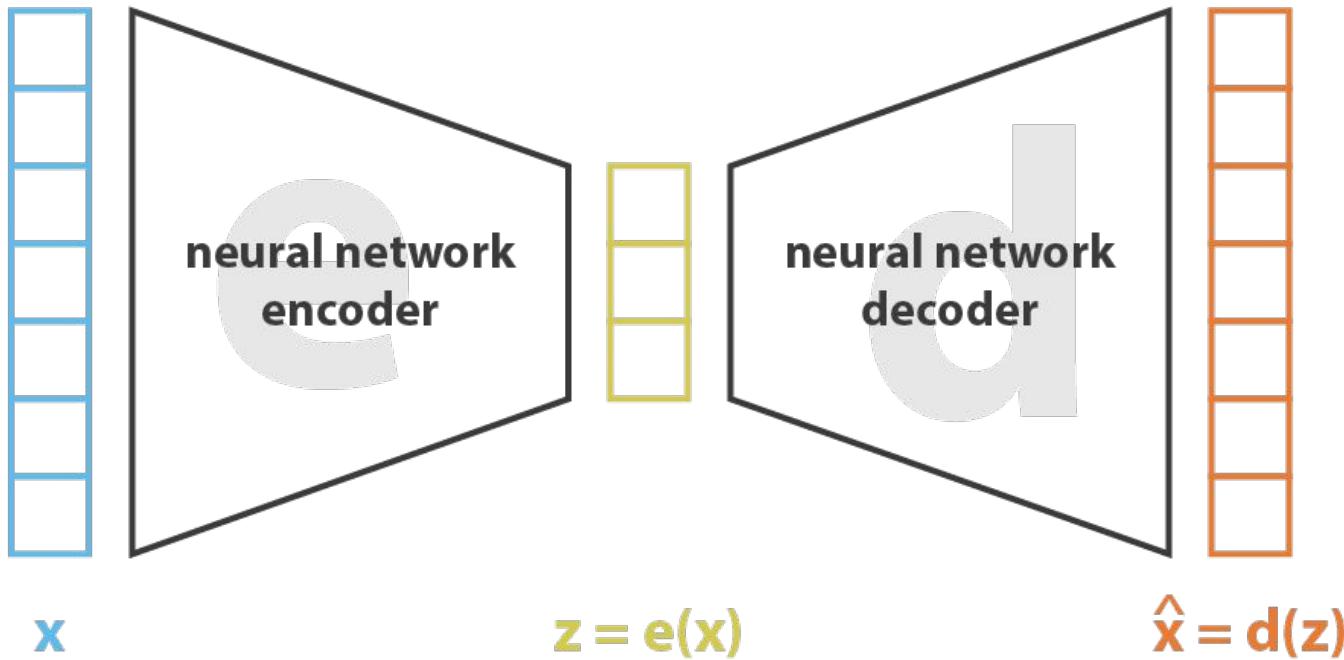


Advantages of Probabilistic PCA

- Interpretation as a ***generative*** model
- **Uncertainty** in latent space
- Can handle **missing data**
- Can "automatically" choose the latent space dimension
- Naturally extendable to **binary/count/categorical data**

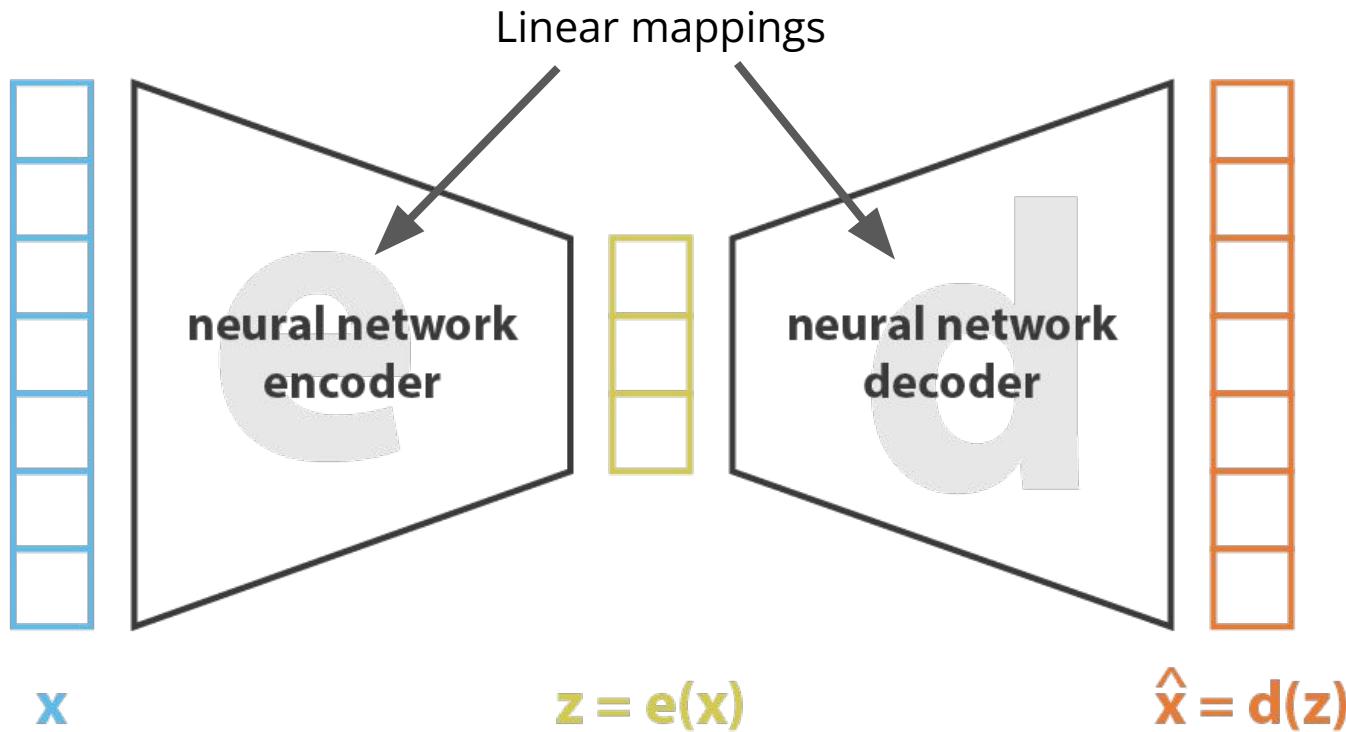
From PCA to Autoencoders

What is an autoencoder?

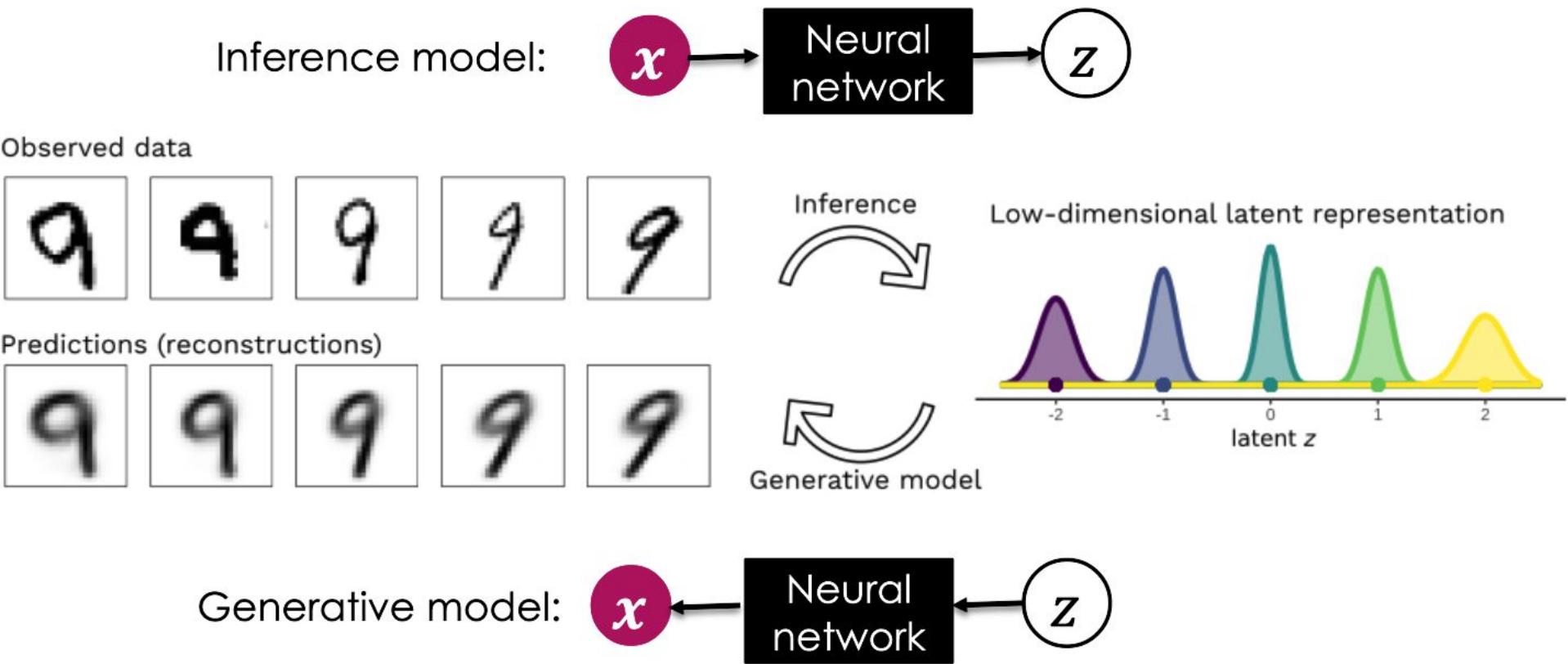


From PCA to Autoencoders

Linear autoencoder implements PCA

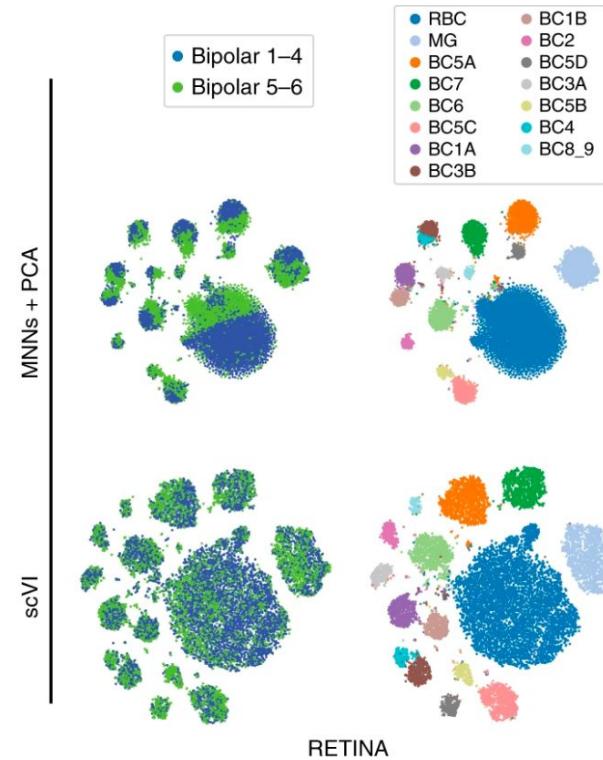
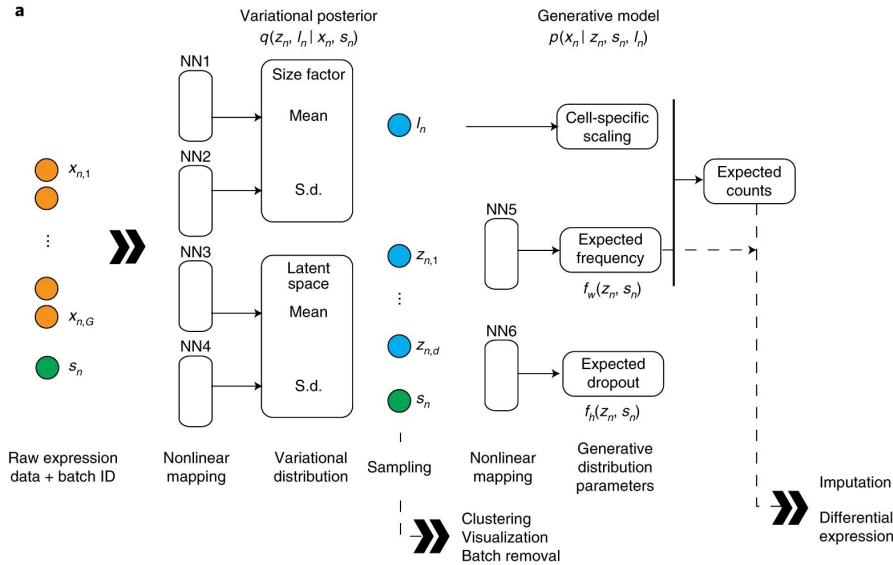


Variational Autoencoders are a non-linear extension of probabilistic PCA



Variational Autoencoders for single cell transcriptomics

Model setup/architecture:



Deep generative modeling for single-cell transcriptomics <https://www.nature.com/articles/s41592-018-0229-2>

scVI-tools

VAE based methods can be useful for large single-cell omics data sets



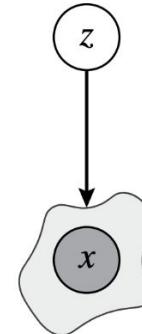
[Get Started](#) [Docs ▾](#) [About ▾](#) [Blog](#) [Discussion ↗](#) [GitHub ↗](#)

Probabilistic models for single-cell omics data

scvi-tools accelerates data analysis and model development, powered by PyTorch and AnnData.

```
pip install scvi-tools
```

[Get Started](#)

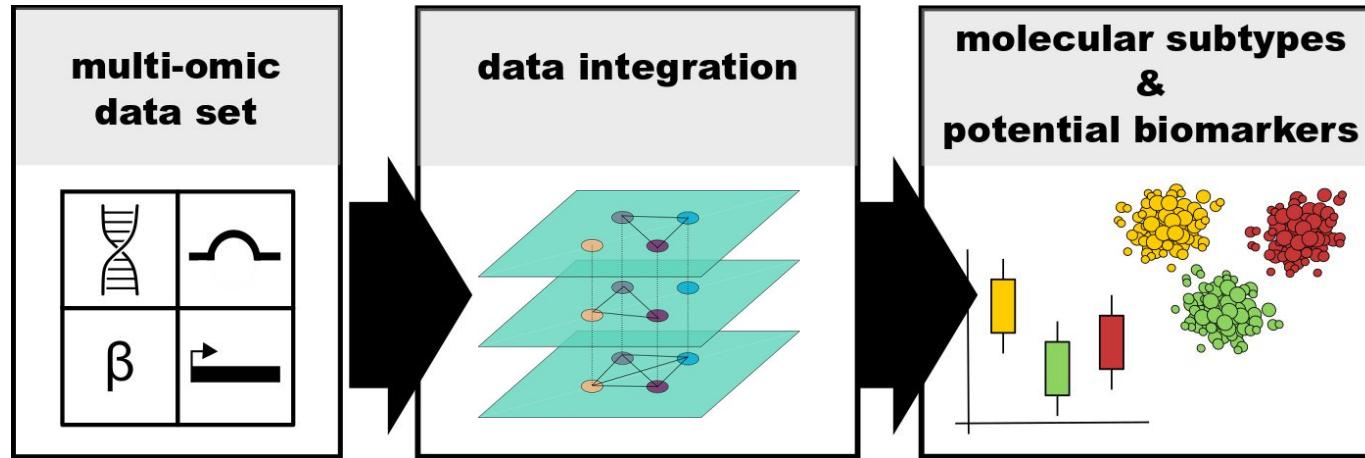


<https://docs.scvi-tools.org/en/stable/tutorials/index.html>

UNSUPERVISED LEARNING

Practical application #3:
Multi-omic data integration

Why are we integrating multi-omic datasets?



What are the challenges of **multi-omic** data integration?



What are the challenges of multi-omic data integration?

Different number
of samples per data-type

Missing values
in feature matrices

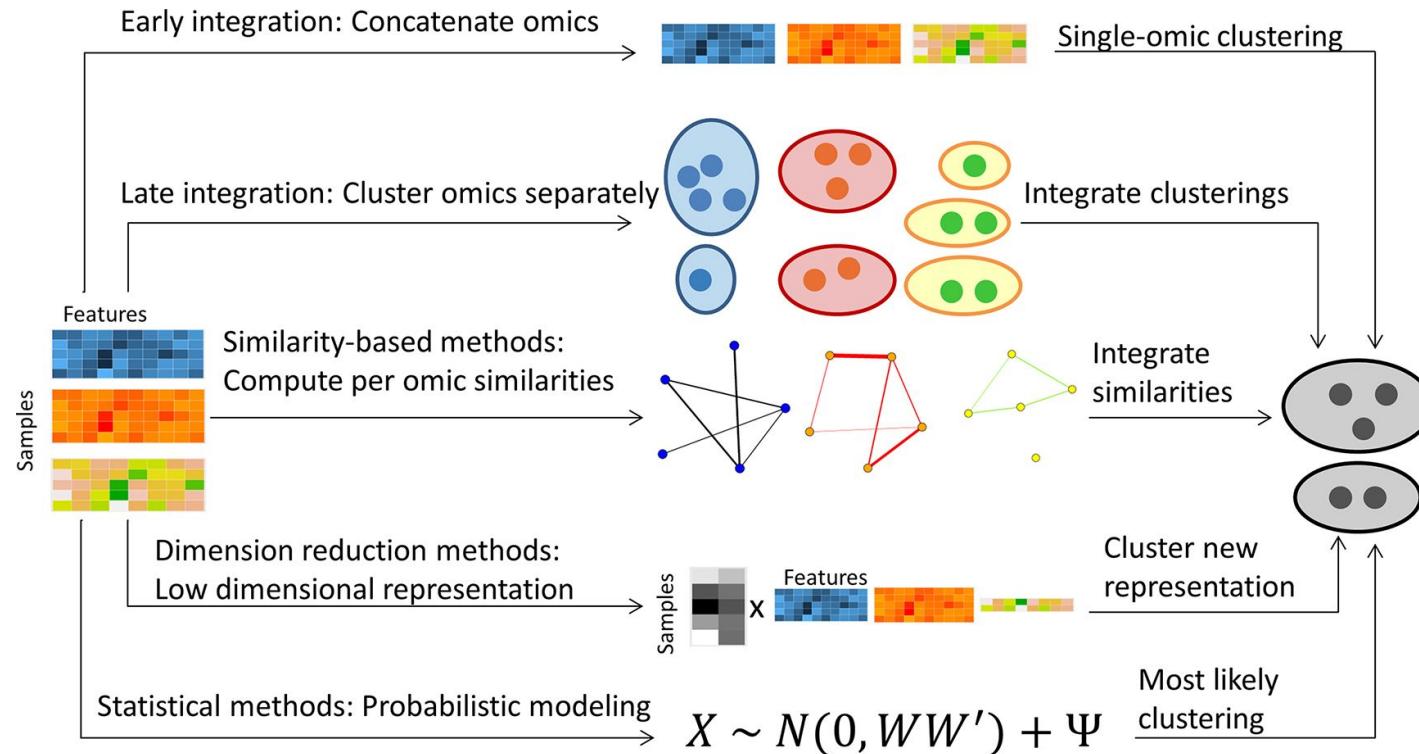
Feature selection

Data heterogeneity

Normalization across
data types following
different distributions

Sparse datasets

Approaches to multi-omic data integration



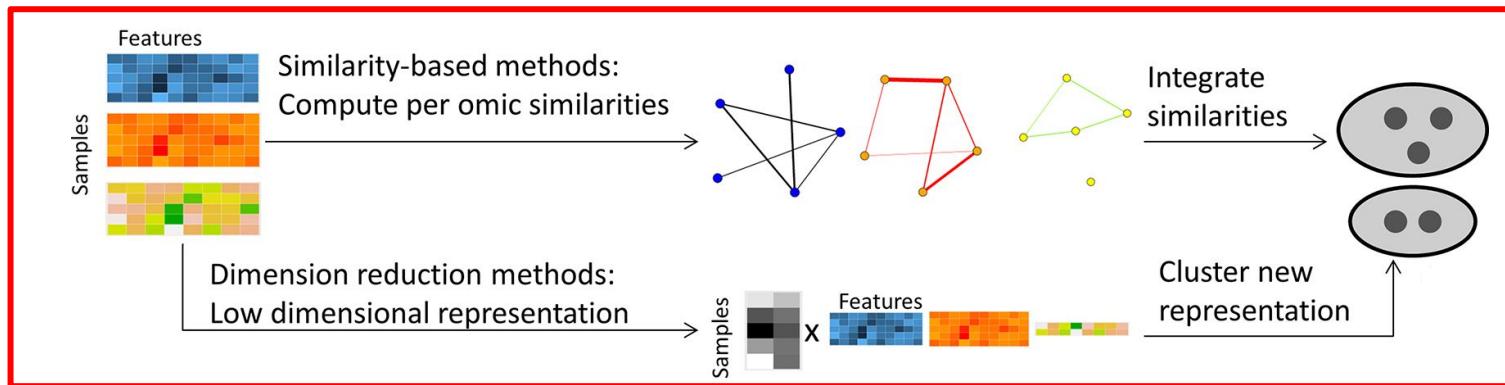
SUMO: Subtyping Using Multi-Omic data



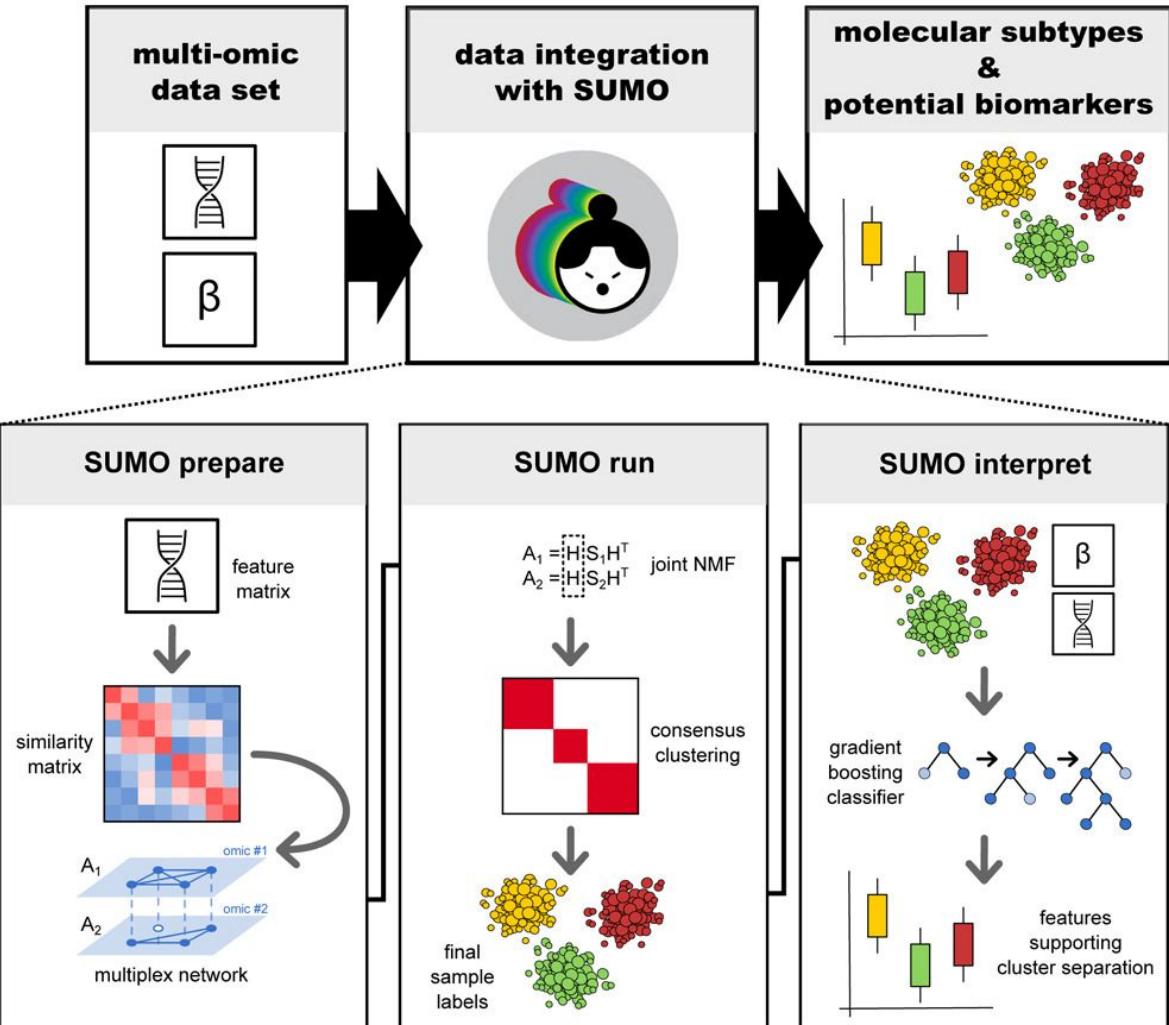
github.com/ratan-lab/sumo



pypi.org/project/python-sumo

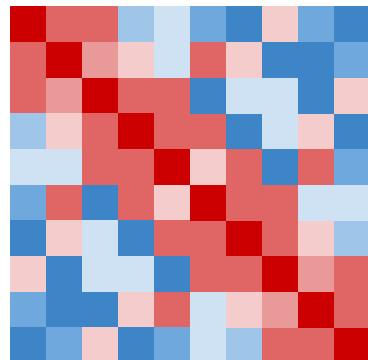
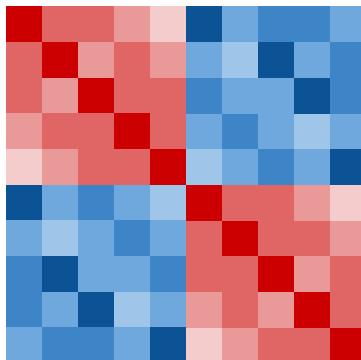
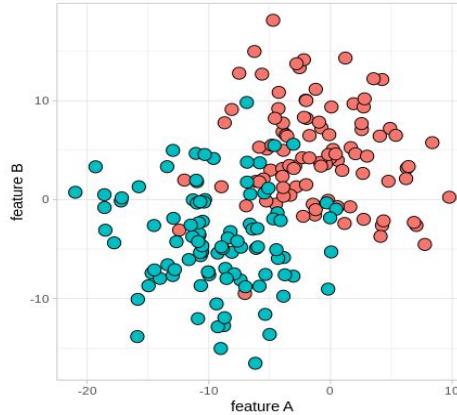
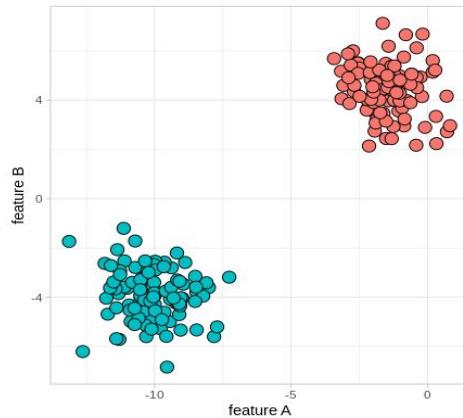


SUMO workflow



Sienkiewicz et al (2022)
Cell Reports Methods

Similarity matrices - intuition



different data type requires
different similarity metrics

SUMO: Factorization

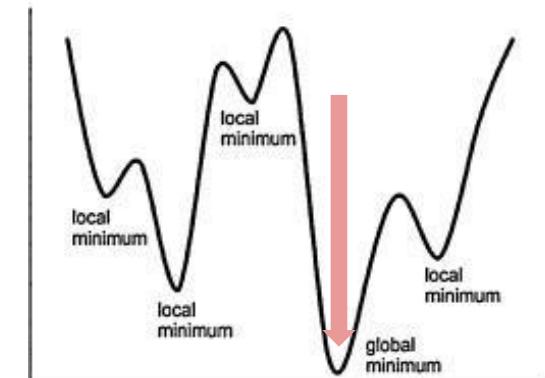
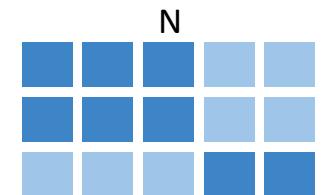
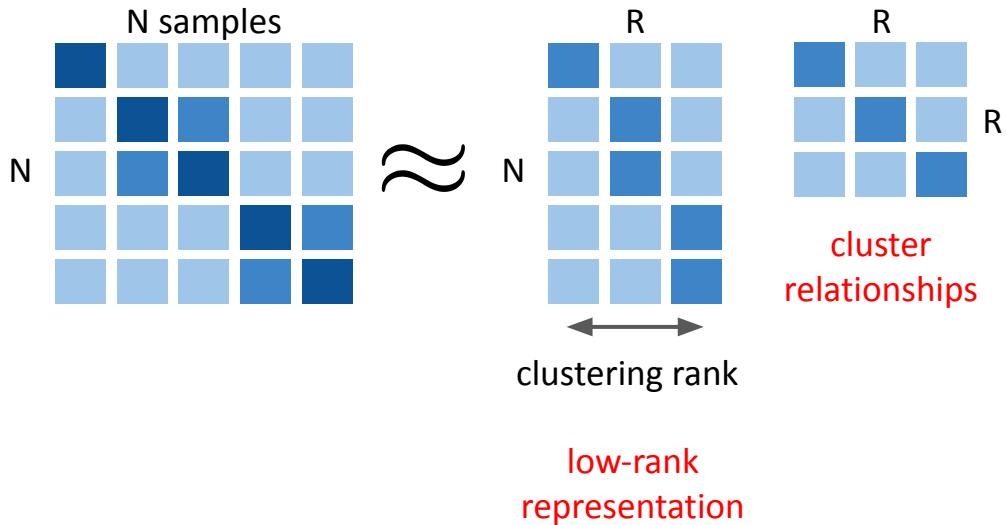
NMF

$$X \approx WH$$

Symmetric
NMF

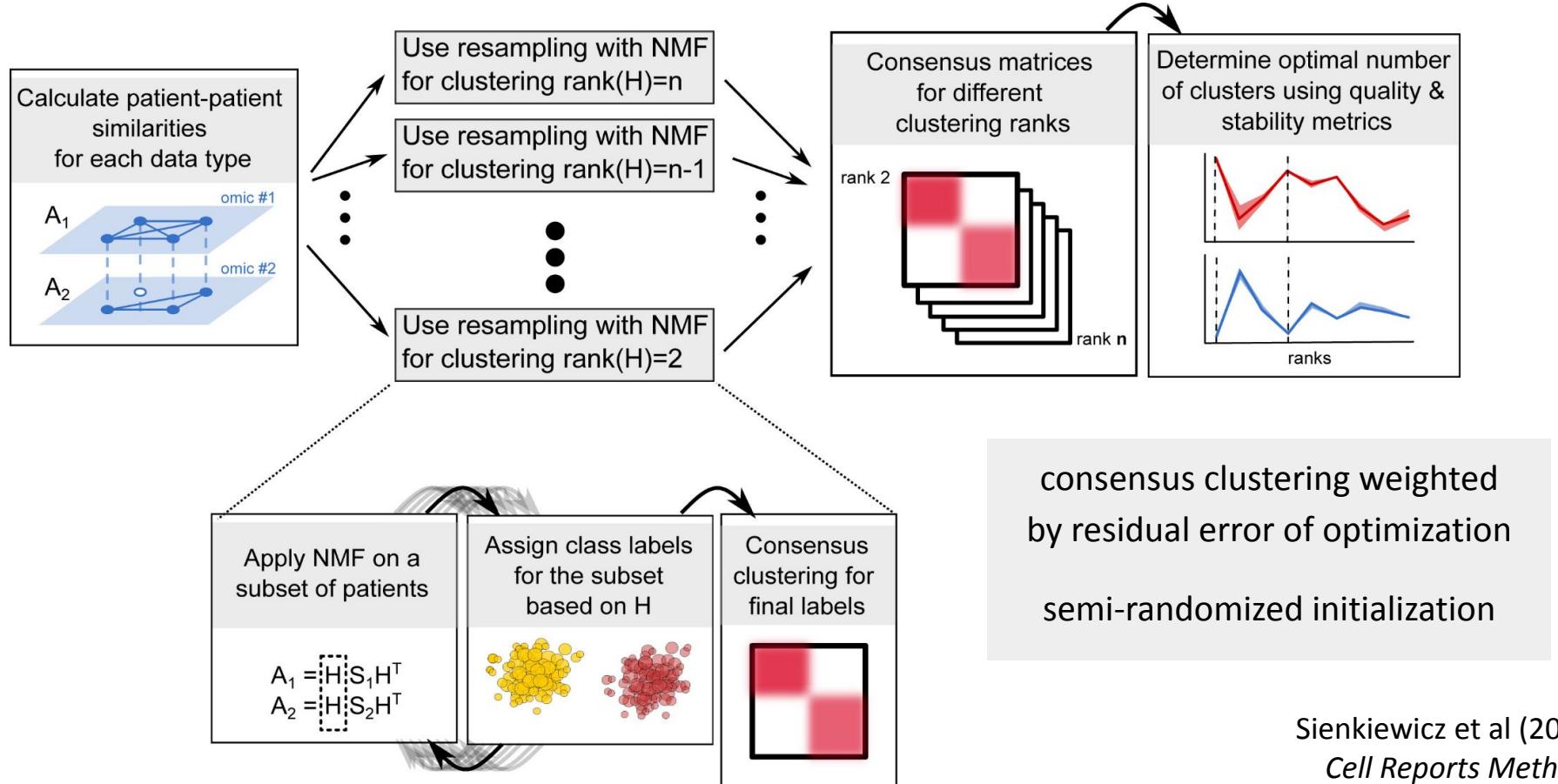
$$A \approx HH^T$$

$$A^i \approx HS_iH^T$$



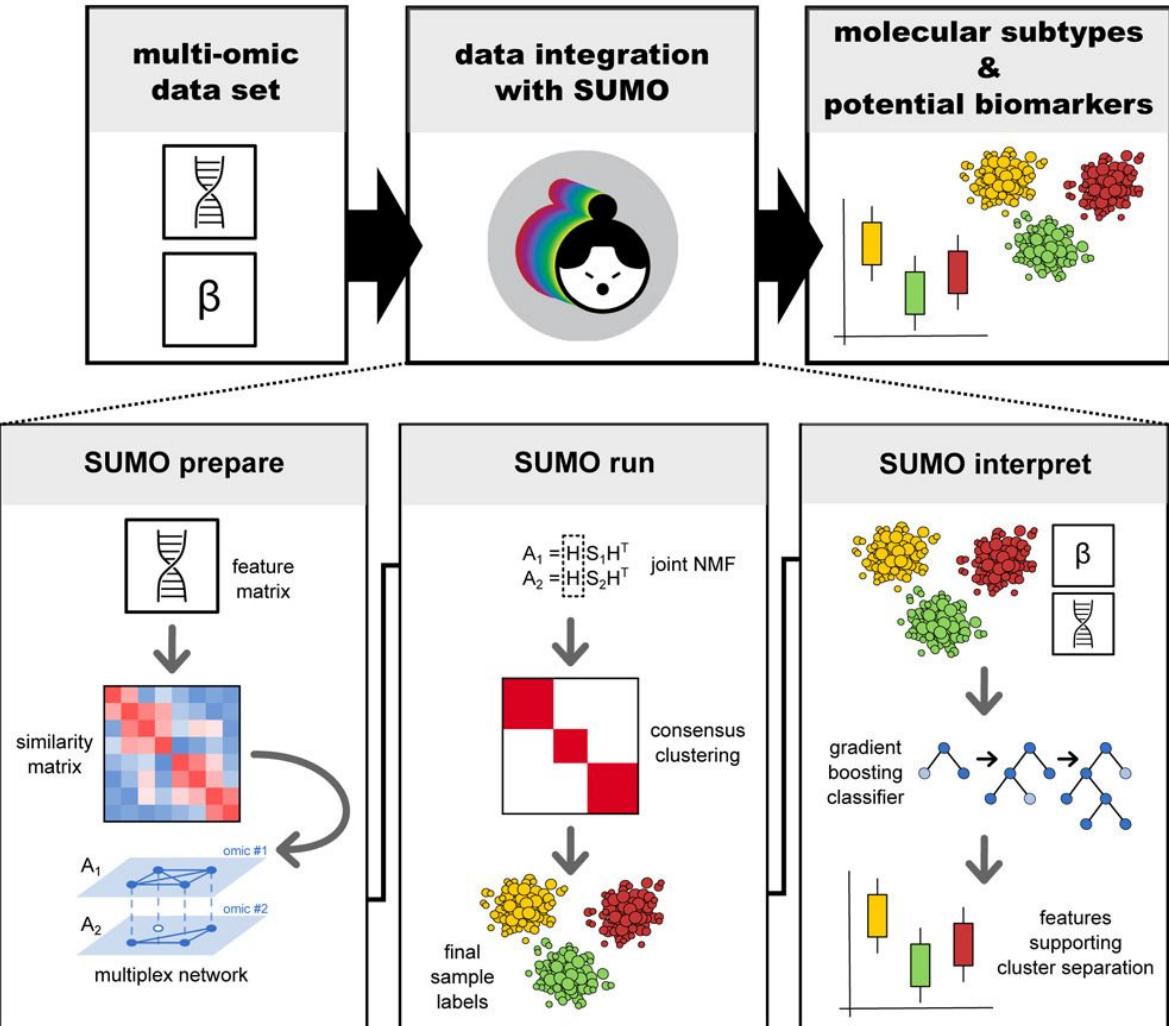
using gradient descent

SUMO: Re-sampling based approach



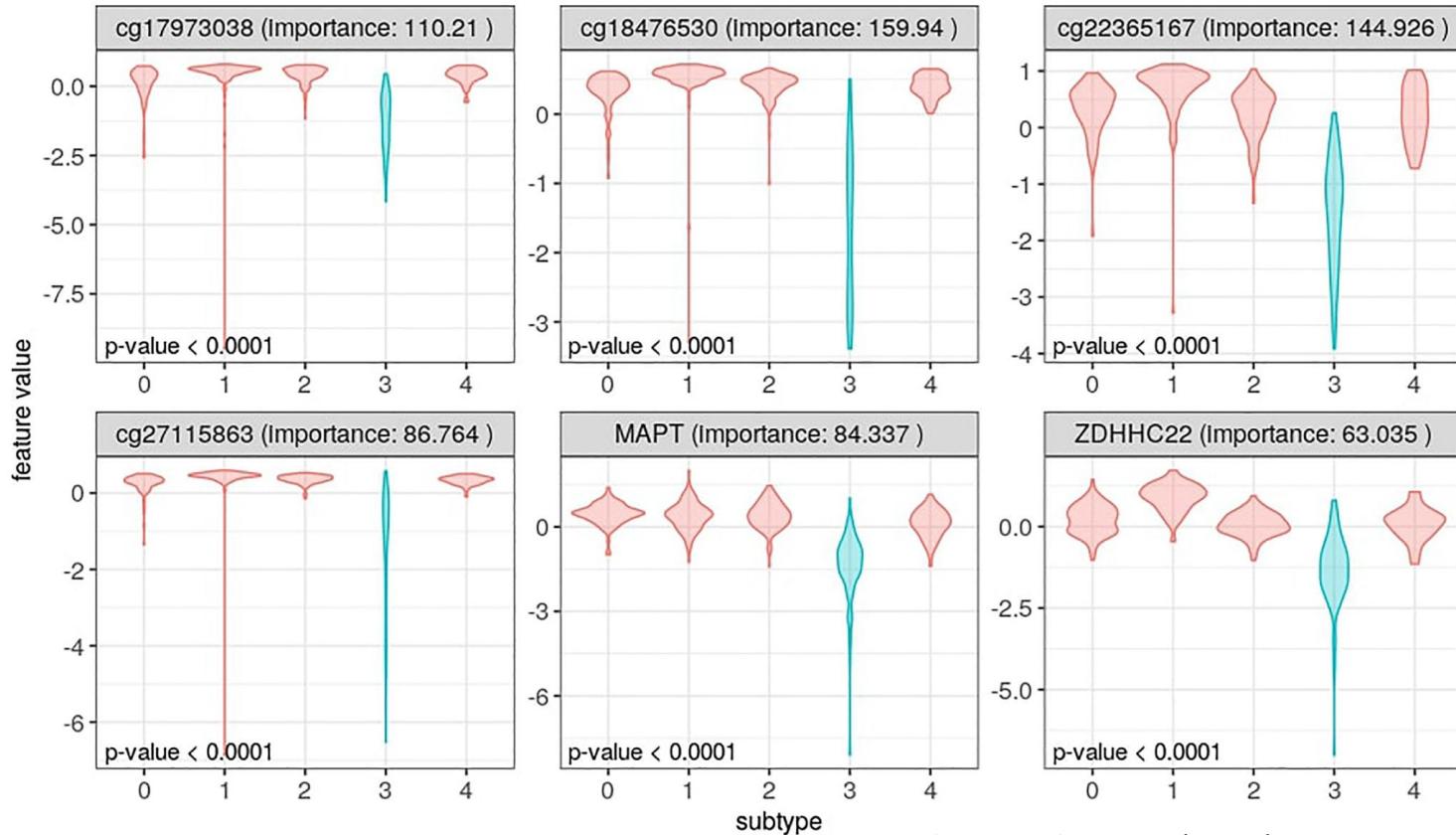
Sienkiewicz et al (2022)
Cell Reports Methods

SUMO workflow



Sienkiewicz et al (2022)
Cell Reports Methods

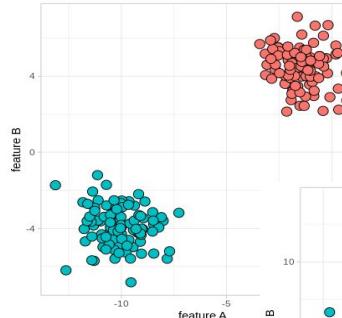
SUMO INTERPRET: downstream analysis



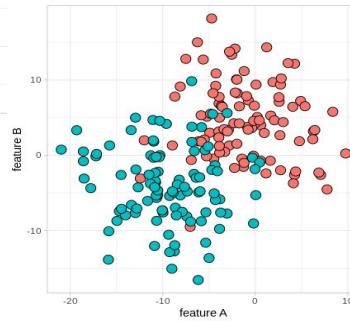
Sienkiewicz & Ratan (2022) STAR Protoc.

How do you **test** a new ML model? - a practical approach

consult your
goals



datasets with
known
ground truth



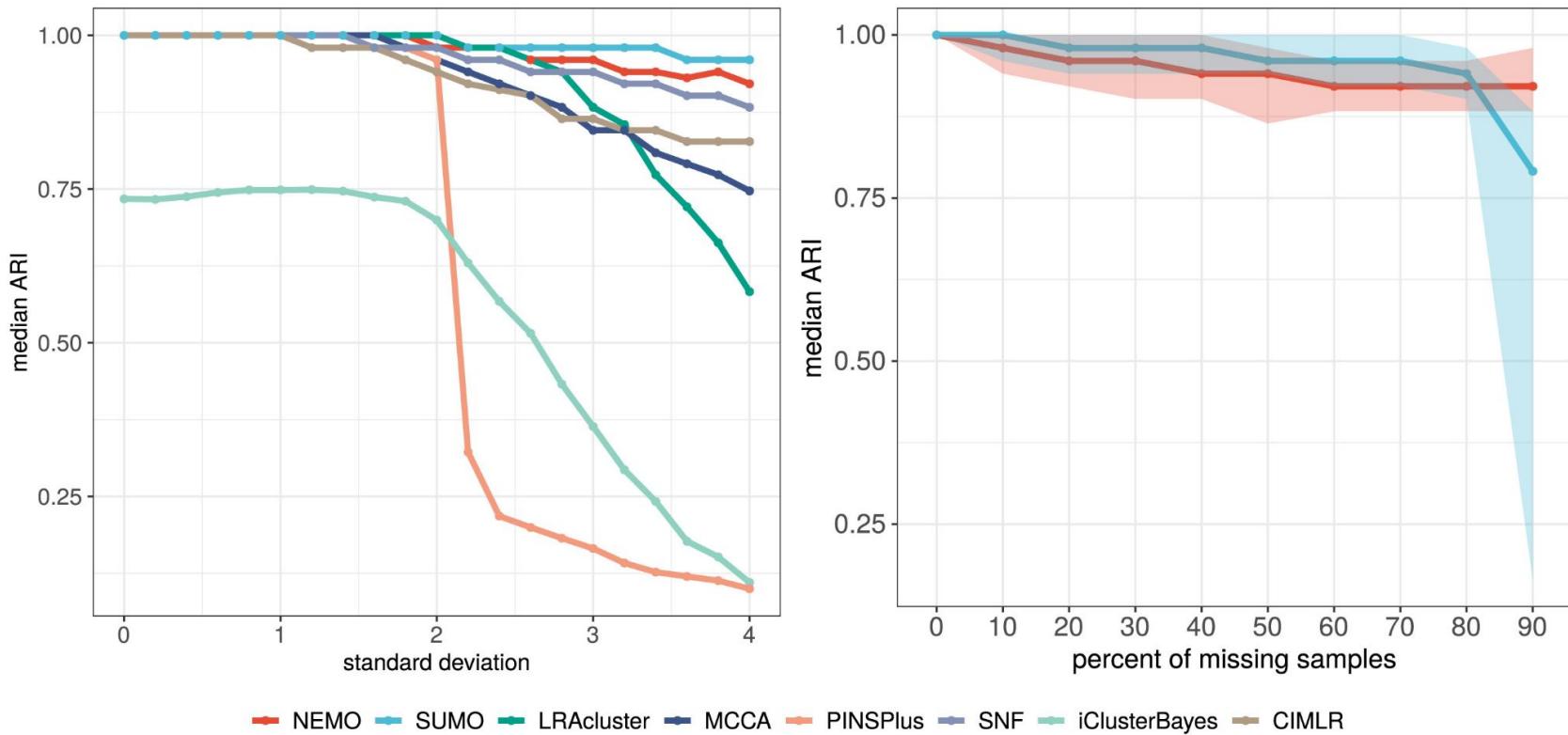
repeat each test



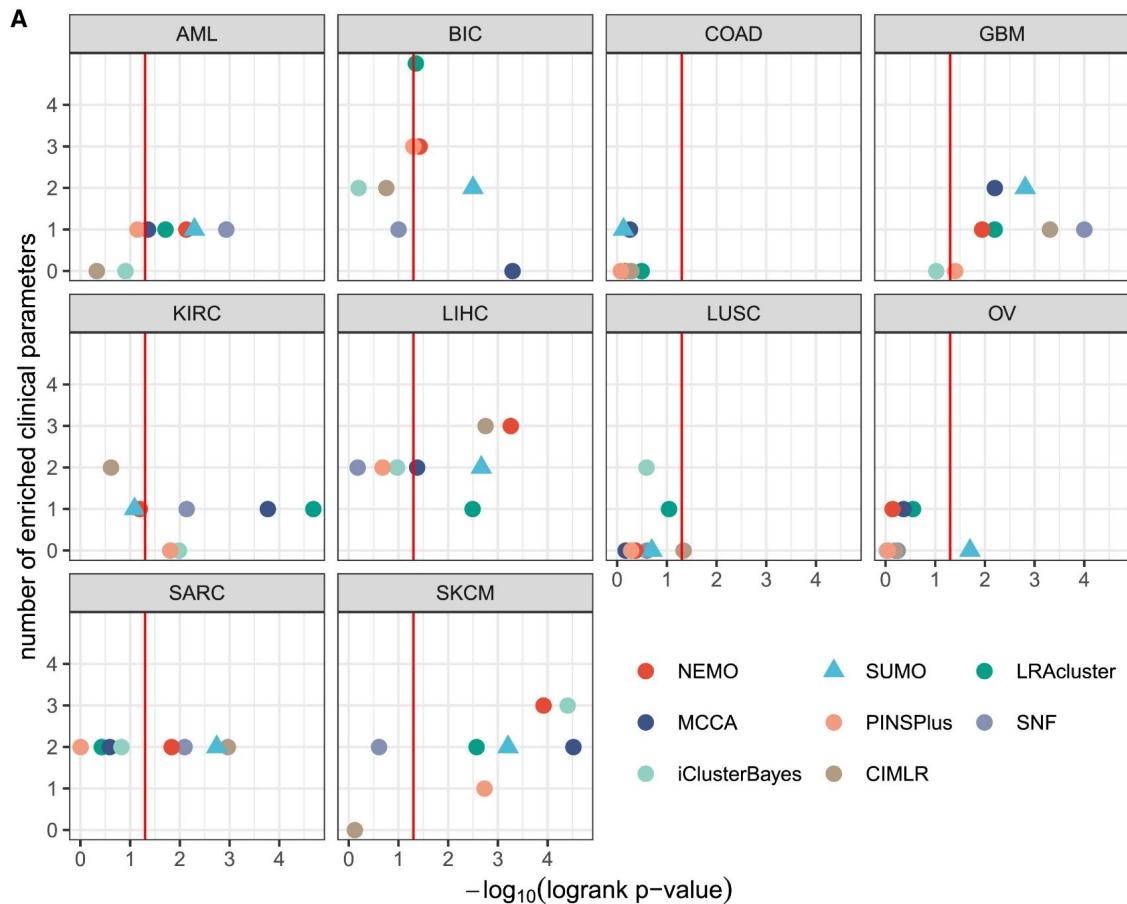
technical considerations
(time & memory)

compare with
competition

SUMO improves performance with noisy and incomplete data

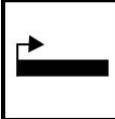
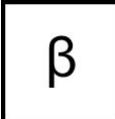
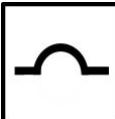


SUMO consistently identifies differential subtypes in TCGA data using an established benchmark

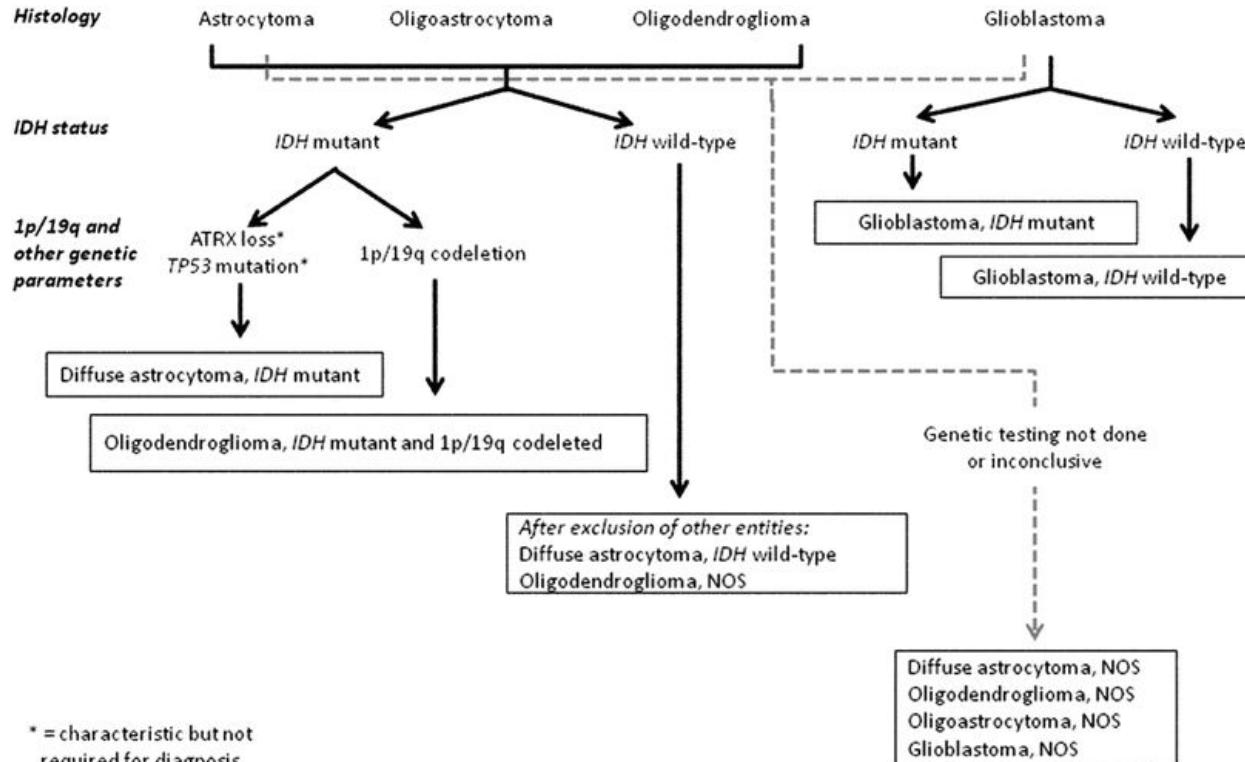


SUMO detects latent relationships between patients in Lower-Grade Gliomas

TCGA-LGG
cohort
(a total of
556 primary tumors)

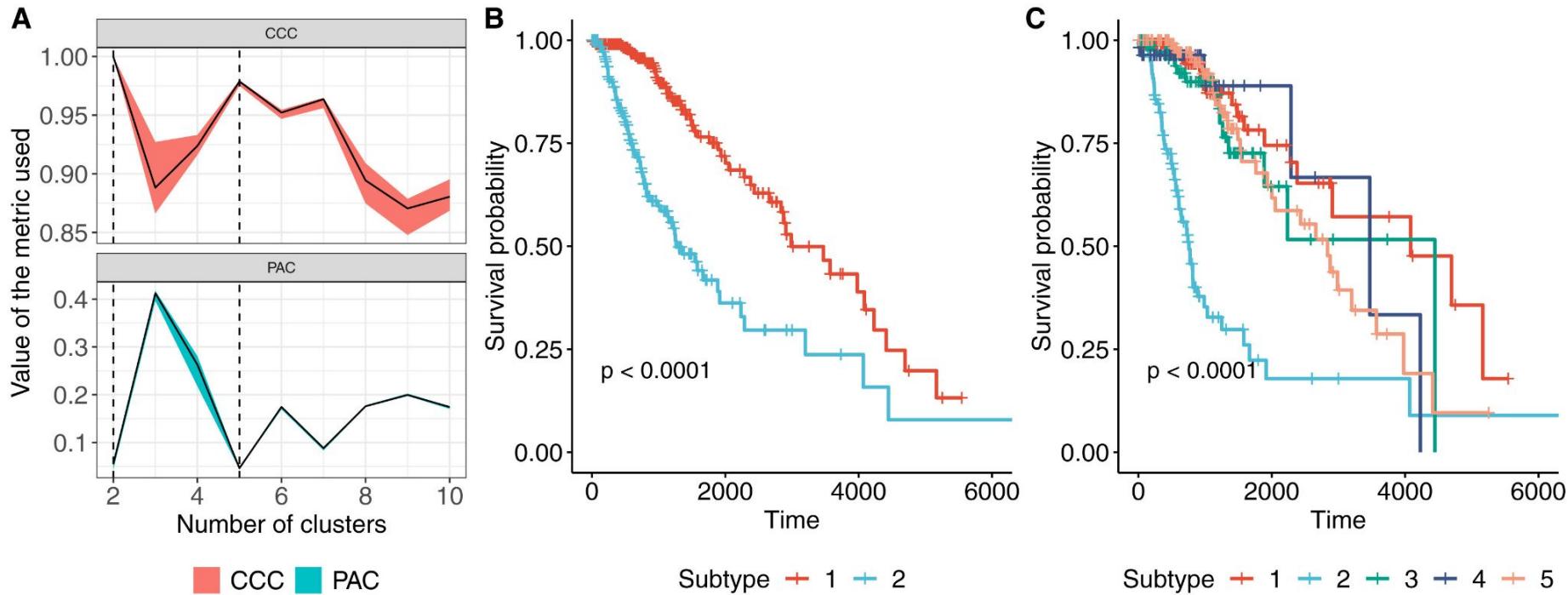
omic	# samples	tech
 RNA-seq	530	Illumina HiSeq
 DNA methylation	530	Illumina Infinium HumanMethylation450
 miRNA expression	524	Illumina HiSeq

WHO guidelines for Low-Grade Gliomas

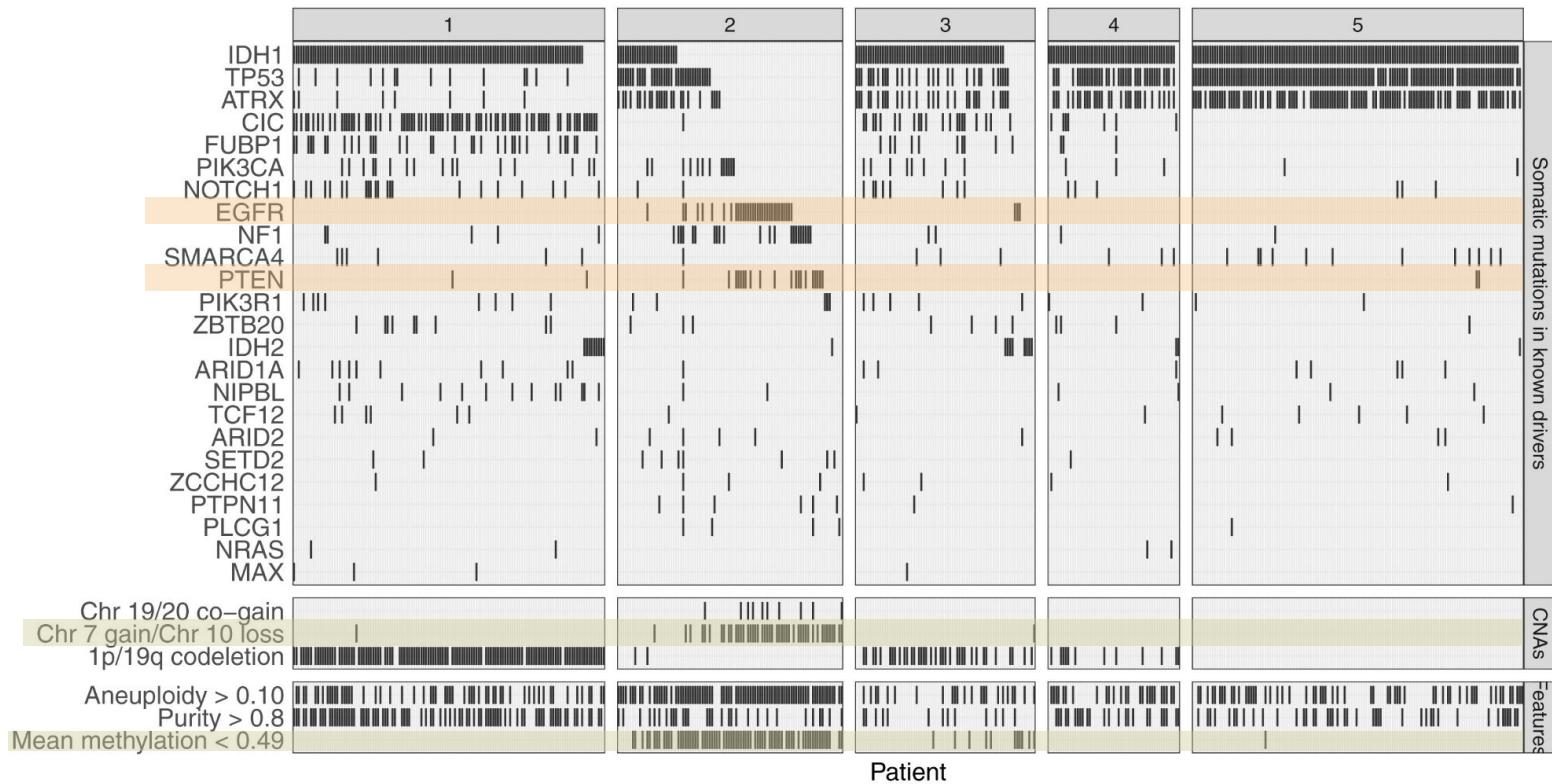


- Gliomas are neuroepithelial tumors originating from glial cells
- Postoperative RT remains the standard of care for adult grade 2/3
- ~80% have mutations in IDH1, 90% of those are R132H.
- ~50% have mutations in TP53, 40% have mutations in ATRX and 20% in CIC

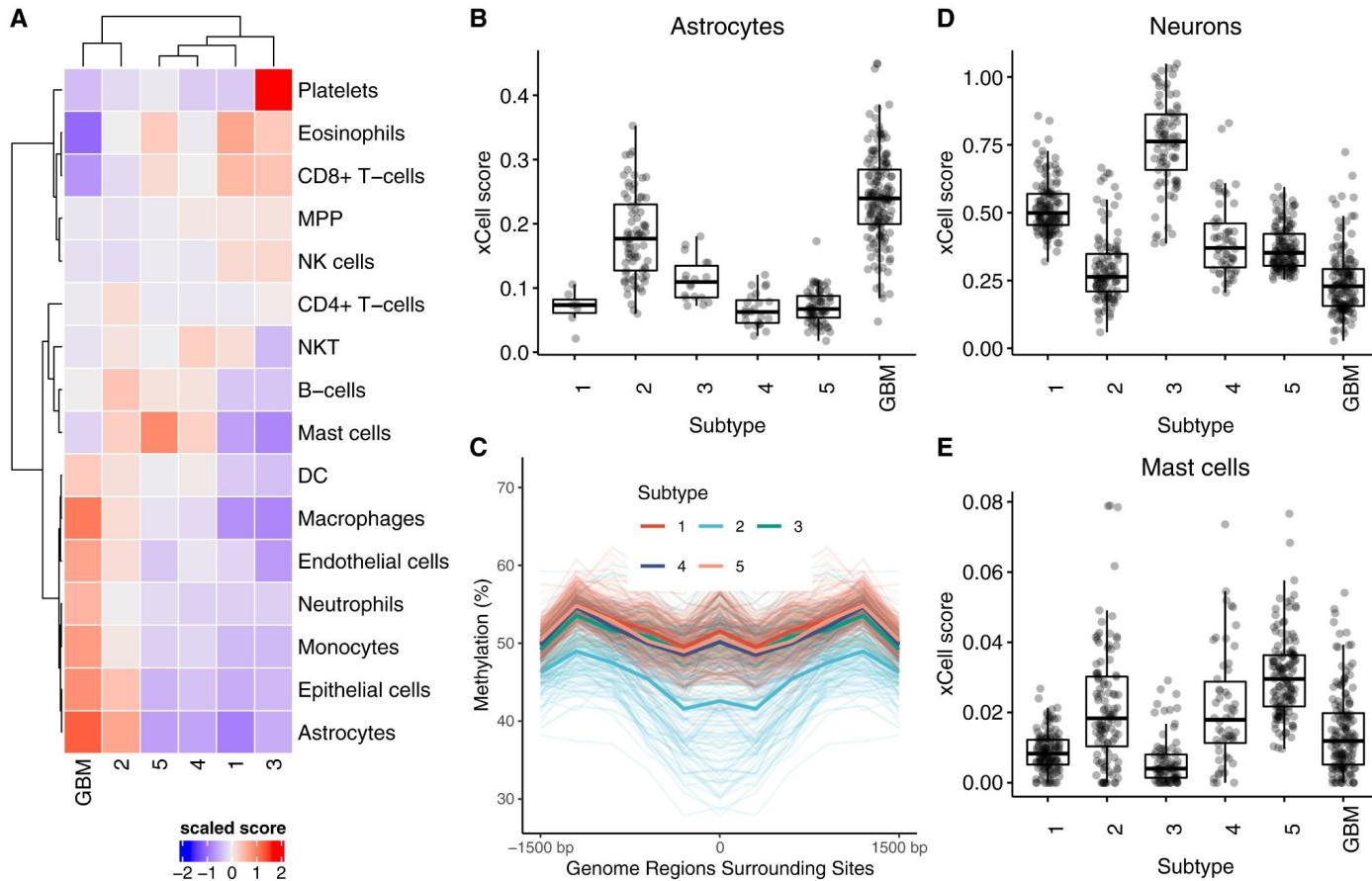
SUMO identifies a subtype of LGG with differential prognosis & GBM-like features



SUMO identifies a subtype of LGG with differential prognosis & GBM-like features

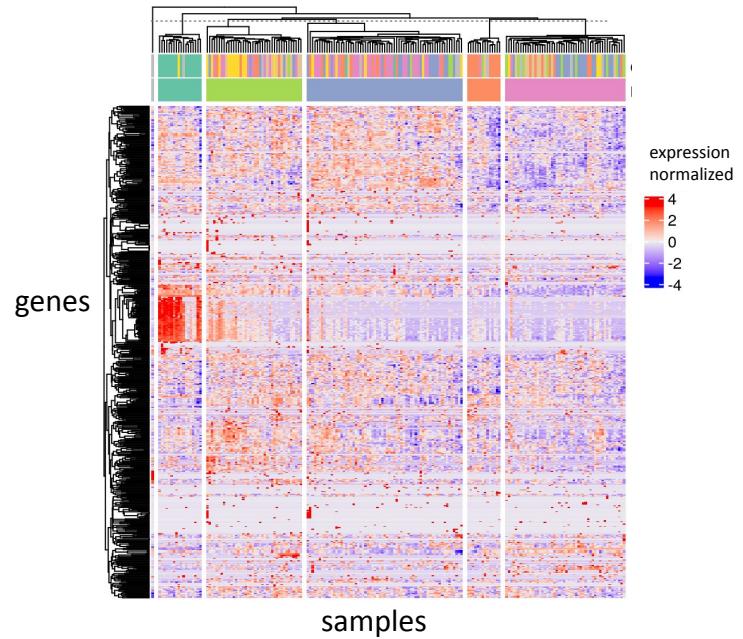


SUMO
subtype 2
cellular profile
is more
similar to
GBMs than
other LGGs



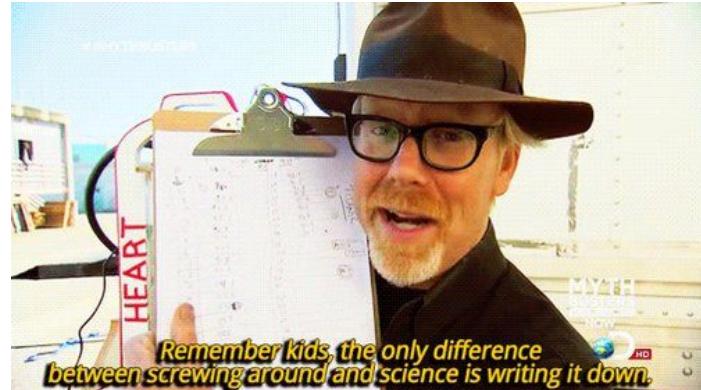
Conclusions

- Heatmaps are your friend



Conclusions

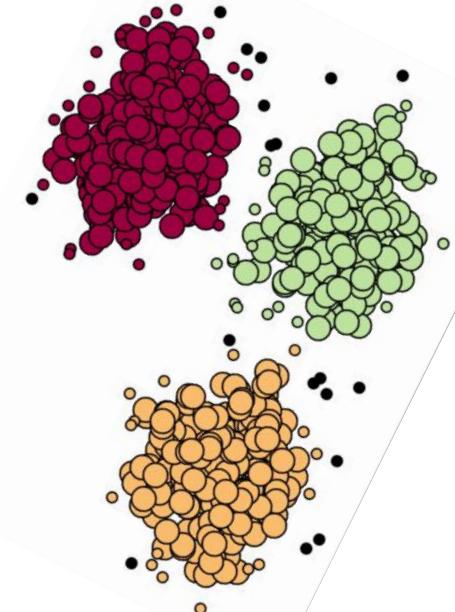
- Heatmaps are your friend
- Analysis-aware experimental design is essential
- Keep track of the metadata



Adam Savage

Conclusions

- Heatmaps are your friend
- Analysis-aware experimental design is essential
- Keep track of the metadata
- Each model has assumptions
- Don't underestimate the value of simulated data



Conclusions

- Heatmaps are your friend
- Analysis-aware experimental design is essential
- Keep track of the metadata
- Each model has assumptions
- Don't underestimate the value of simulated data
- If you are interested in multi-omic data integration, come to my hackathon!

20.09	21.09	22.09
TUE	WED	THU
DAY 5	DAY 6	DAY 7
Breakfast	Breakfast	Breakfast
Feature selection Miron Kursa	Deep Learning in practice Dmytro Fishman	Deep Learning in practice Dmytro Fishman
Coffee break	Coffee break	Coffee break
Survival analysis German Demidov	Deep Learning In practice Dmytro Fishman	Deep Learning In practice Dmytro Fishman
Lunch	Lunch	Lunch
	Networking / City Sightseeing	
Coffee break		Coffee break
		
Dinner		Dinner

Acknowledgements

SUMO TEAM

University of Virginia

Ajay Chatrath

John Lawson

Nathan Sheffield

Aakrosh Ratan

National University of Singapore

Jinyu Chen

Louxin Zhang

Project funding



National Institutes
of Health



NIGMS



Tri-Institutional PhD Program Computational Biology & Medicine



Memorial Sloan Kettering
Cancer Center



Weill Cornell
Medicine

The Mason Lab