

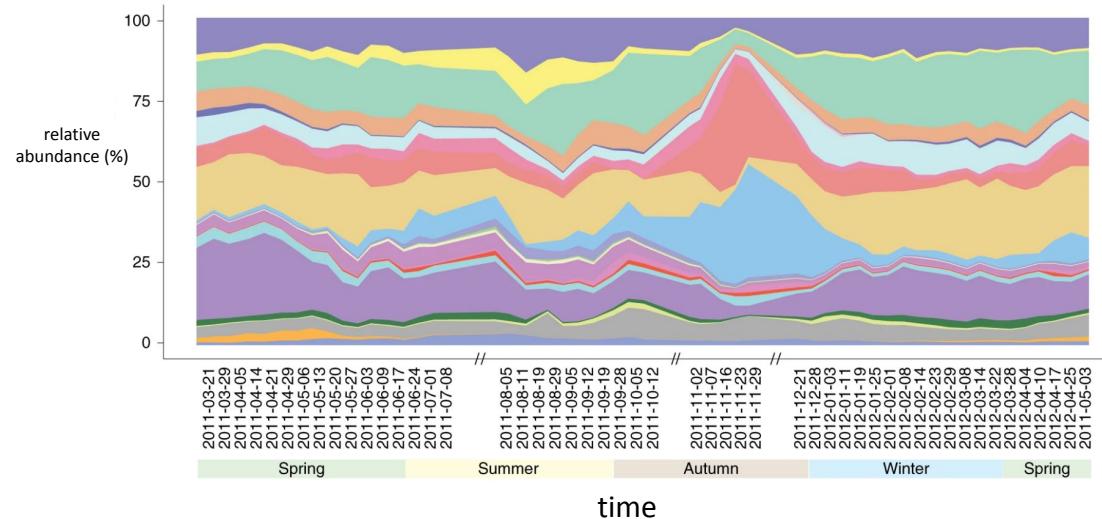
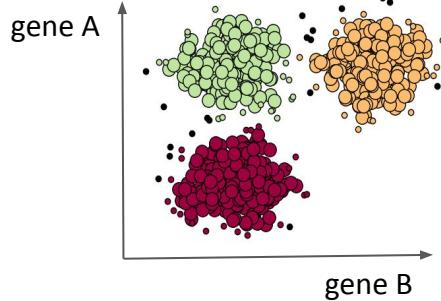
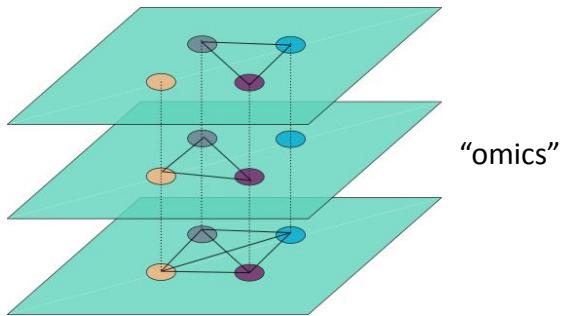
# Multi-omics sequencing data integration *(in practice)*

**Karolina Sienkiewicz**

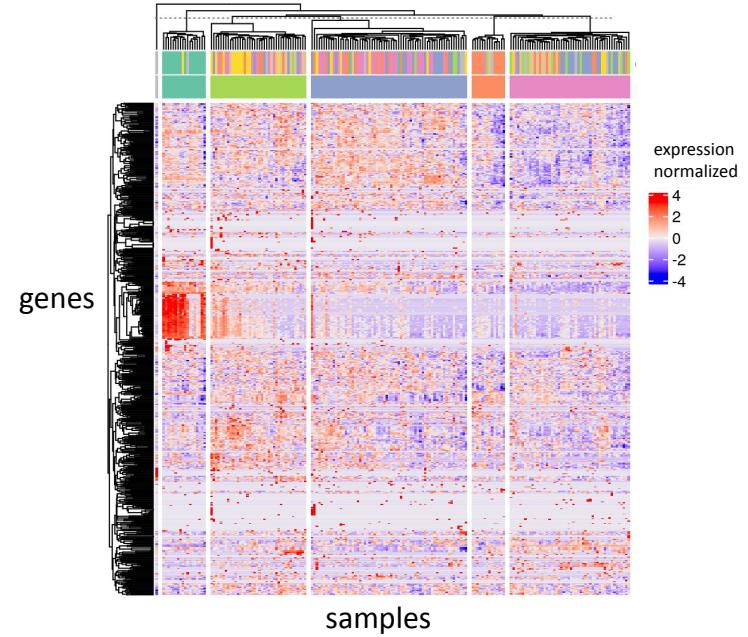
Mason Lab & Melnick Lab  
Weill Cornell Medicine

Single molecule sequencing class  
14th May, 2024

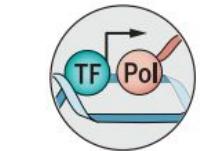
# Biological datasets are complex



PMID:  
33139880

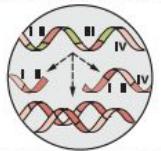


# Single-cell epigenetic profiling

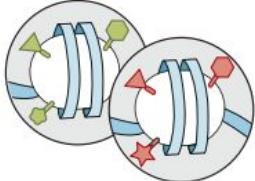


## Transcription factor binding

TF binding interacts with DNA methylation and chromatin accessibility



## Transcription and RNA maturation



## Histone modifications

Modifications can be active marks (e.g., H3K4me3 in green) or repressive marks (e.g., H2K27m3 in red)



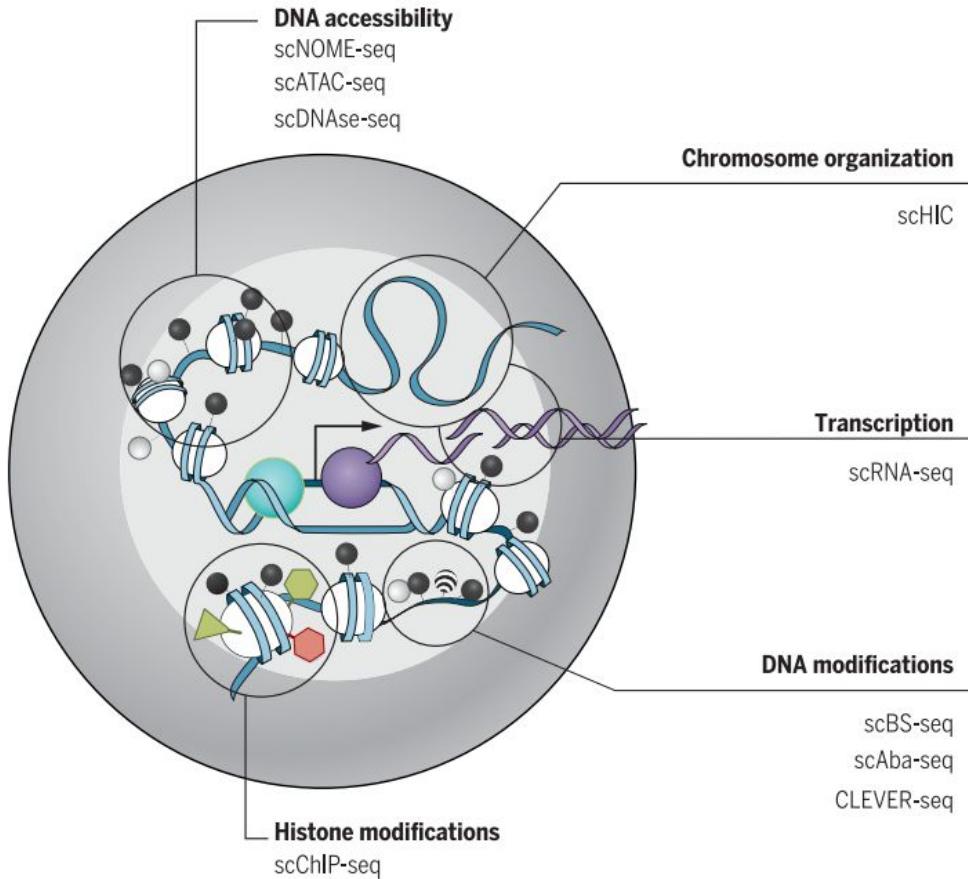
## DNA modifications

● C      ● 5mC  
● 5hmC / 5fC / 5caC

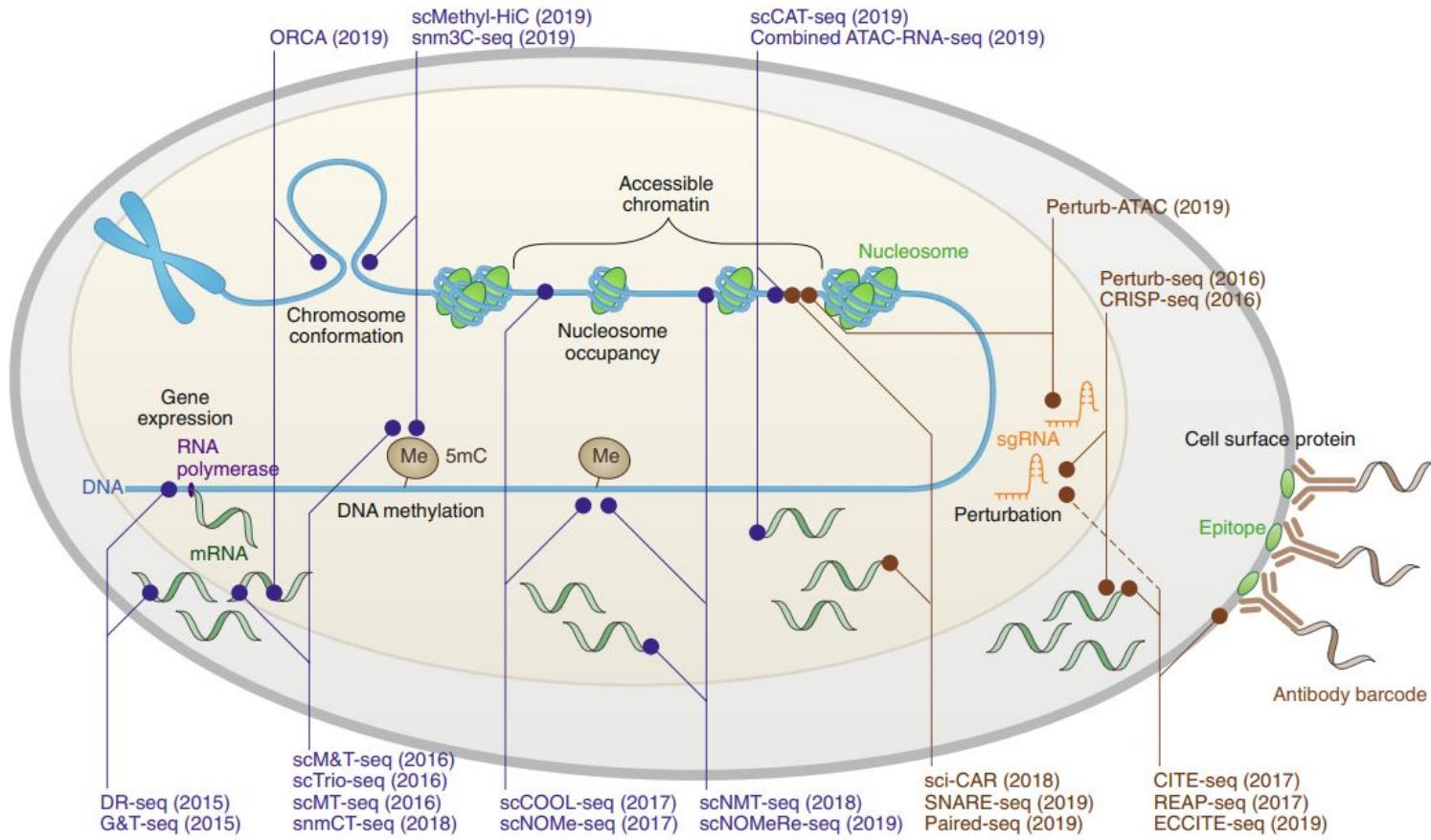


## Chromosome organization

Higher-order chromatin organization into LADs and TADs



# Single-cell multimodal analysis



# 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

Data heterogeneity

Missing values  
in feature matrices

Batch effect

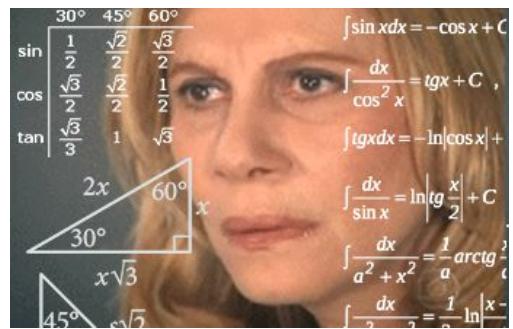
Normalization across  
data types following  
different distributions

Feature selection

Sparse datasets

# Batch effect = non-biological variation

Why is there  
a batch effect  
in my data?

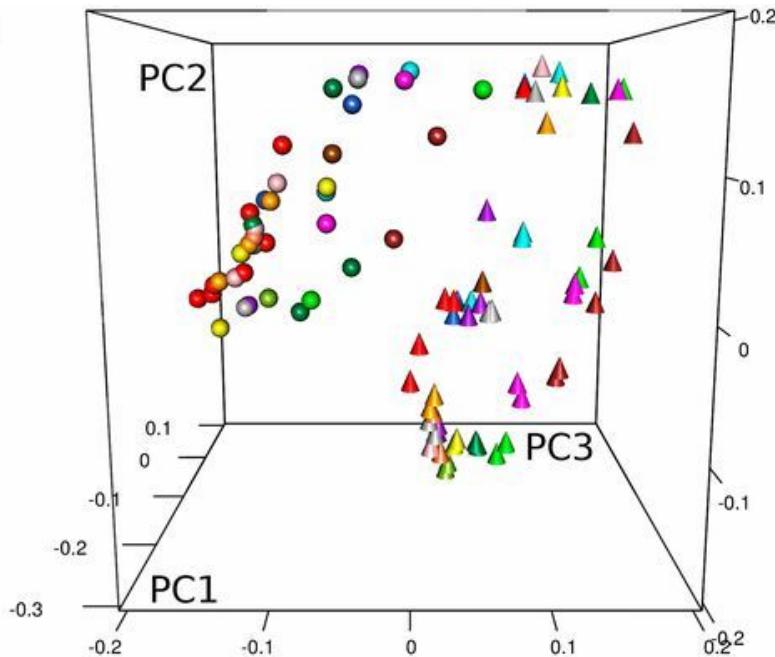


- Different days/months of the experiments
- Different reagents (enzymes, buffers)
- Different mice (from different companies)
- Different sequencers
- Lab protocol or experimenter/technicians

# Batch effect example

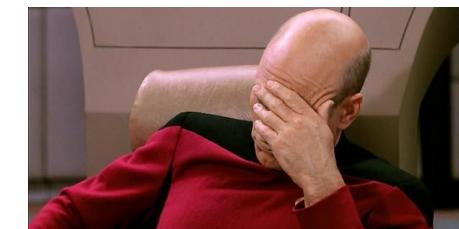
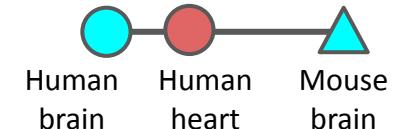
“Gene Expression Is More Similar Among Tissues Within a Species Than Between Corresponding Tissues of the Two Species”

A



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)

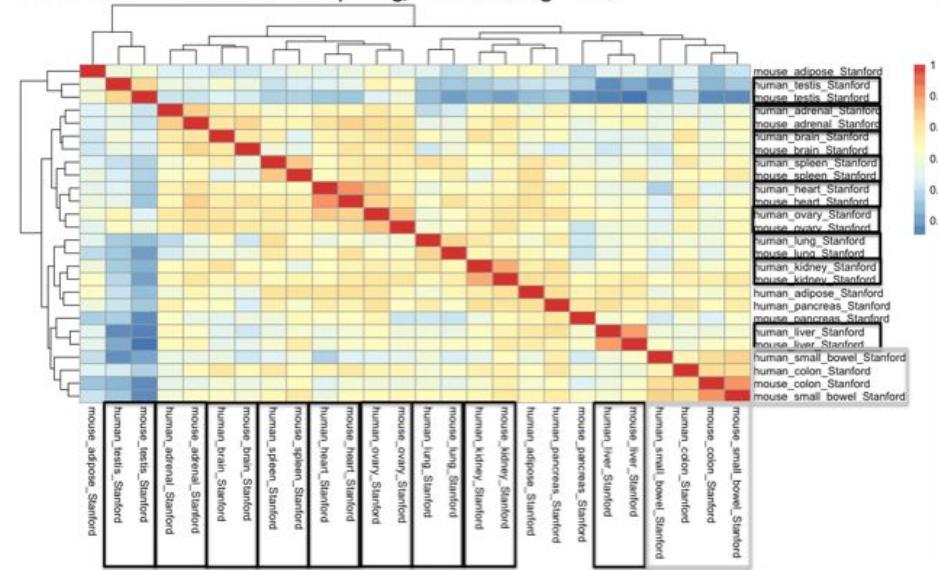
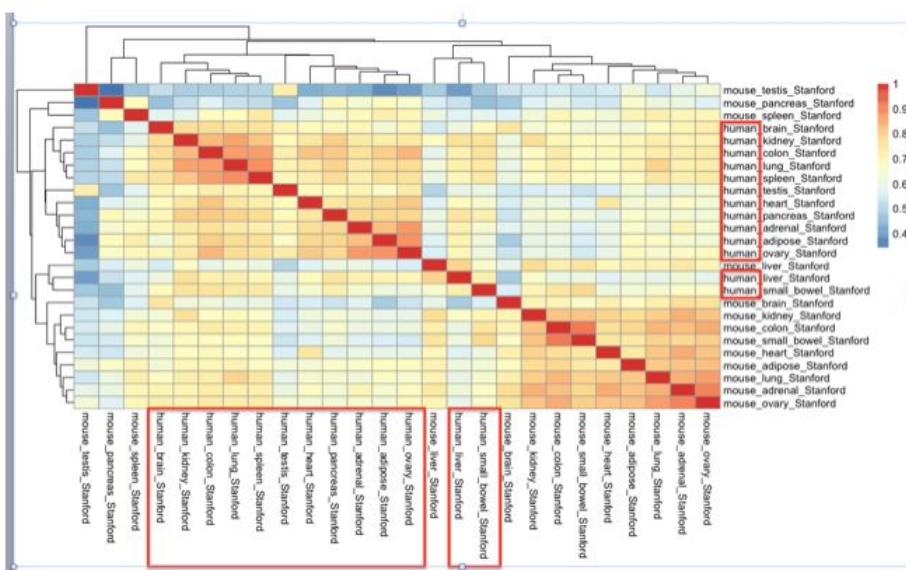




Yoav Gilad  
@Y\_Gilad

We reanalyzed the data from  
[pnas.org/content/111/48...](https://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

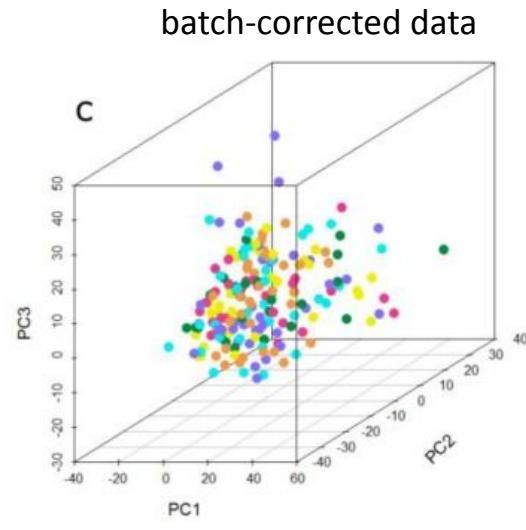
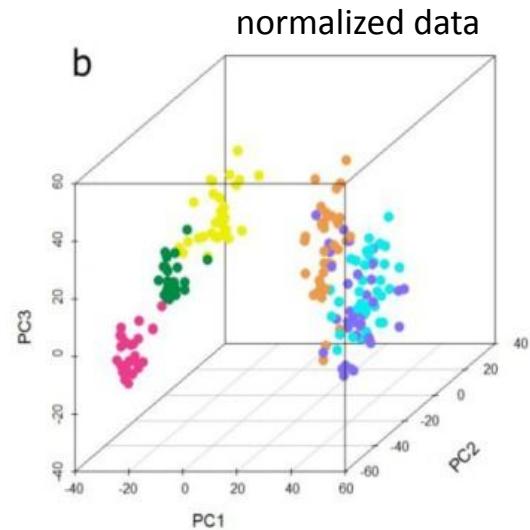
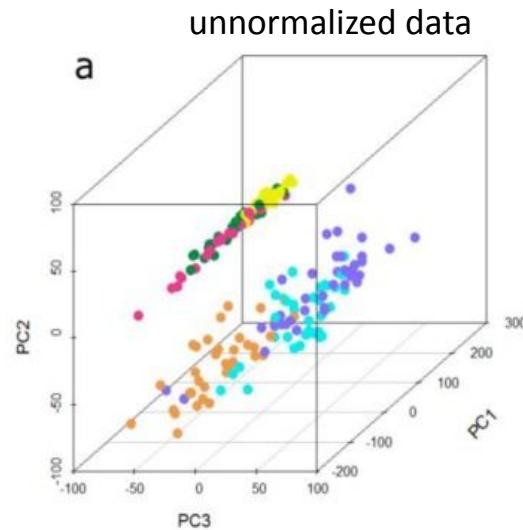
balancing  
groups of interest

avoiding  
“perfect confounding”



# Combating batch effect

visualize/cluster → identify → correct



**ComBat**

Johnson et al (2007) *Biostatistics*

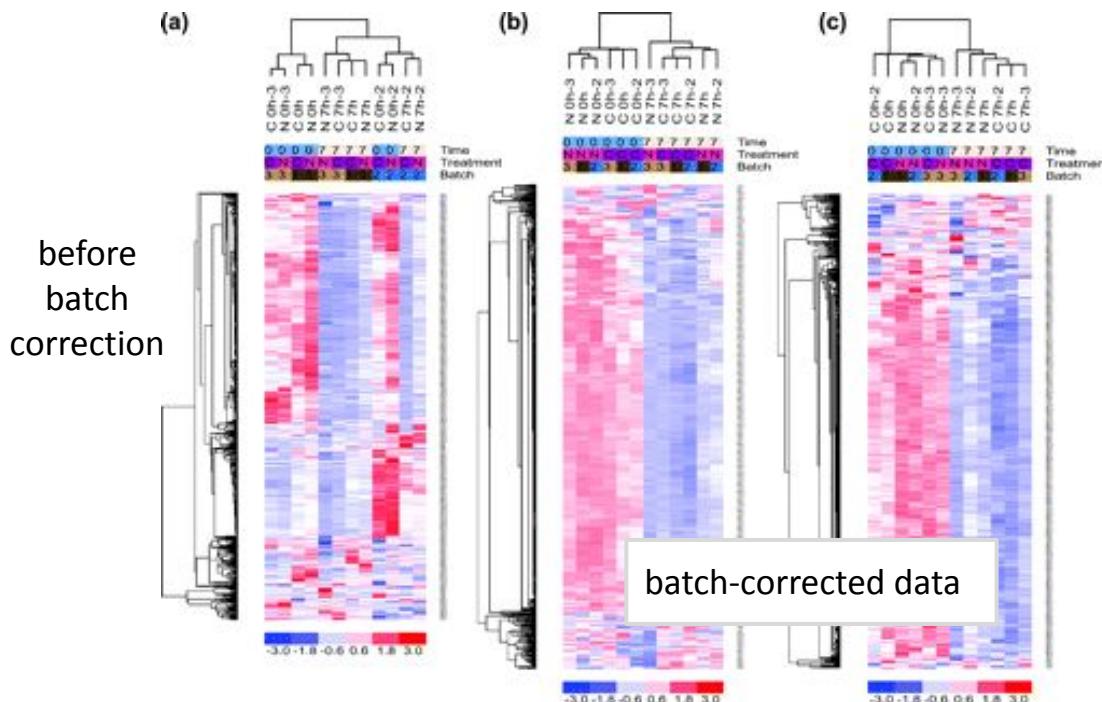
**limma**

complex (e.g. nested) batches

Brezina et al (2015) *Microarray*

# Combating batch effect

visualize/cluster → identify → correct



## ComBat

Johnson et al (2007) *Biostatistics*

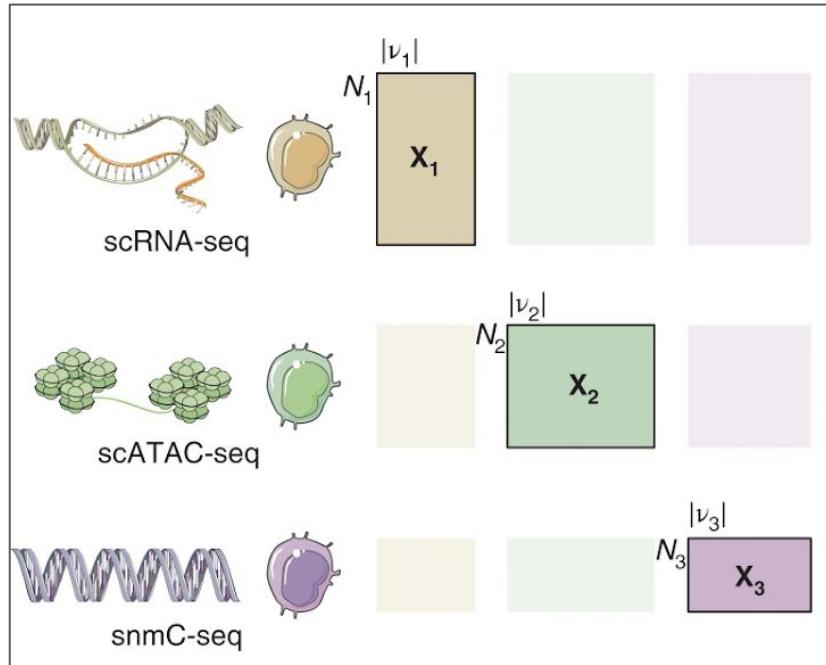
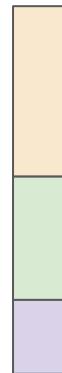
## limma

complex (e.g. nested) batches

Johnson et al (2007) *Biostatistics*

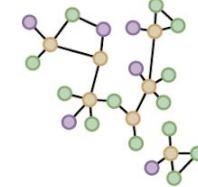
# Major obstacles of integrative modeling: distinct feature/sample space

conversion to  
common feature  
space

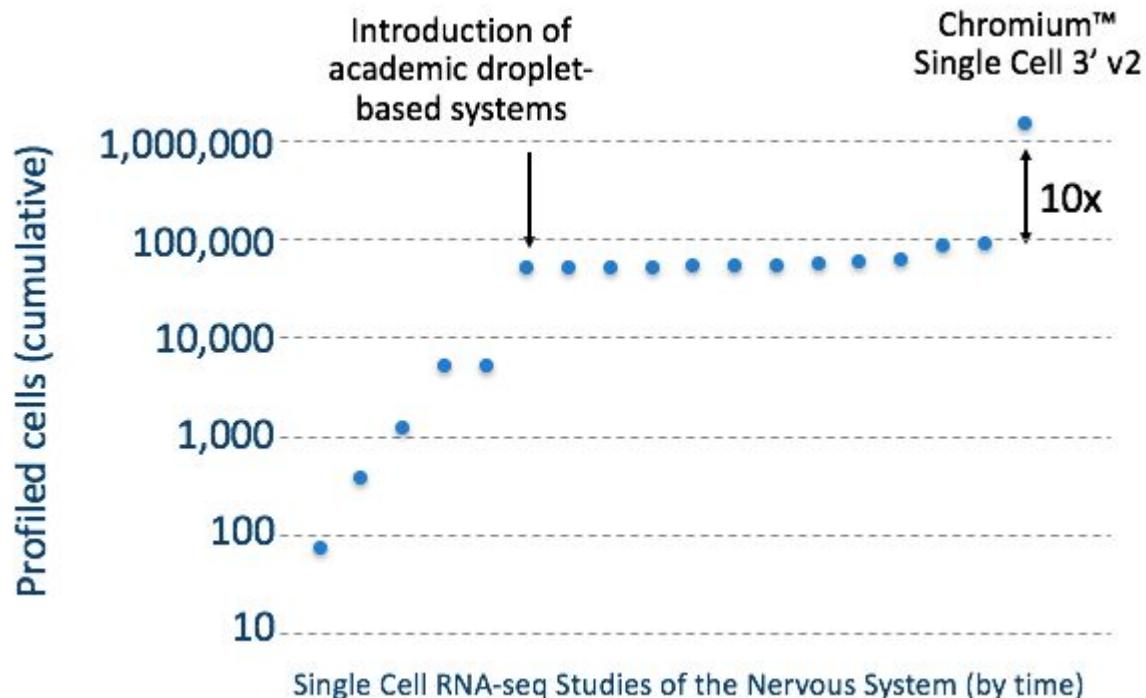


joint matrix factorization

cell matching

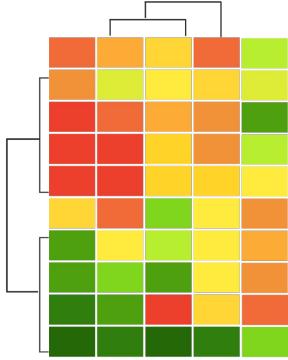


## Major obstacles of integrative modeling: data scale



Approximate cumulative number of cells profiled in scRNA-seq literature (up to Oct 2016). *10x genomics*

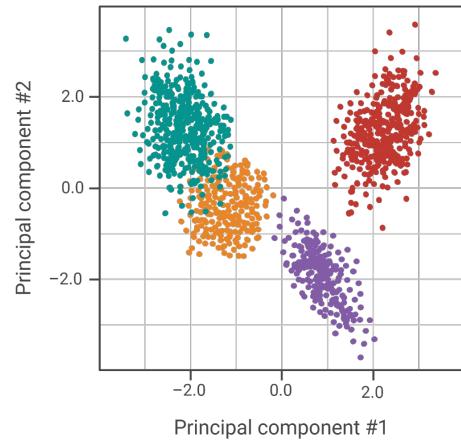
# Primer on dimensionality reduction



data  
representation  
(full)



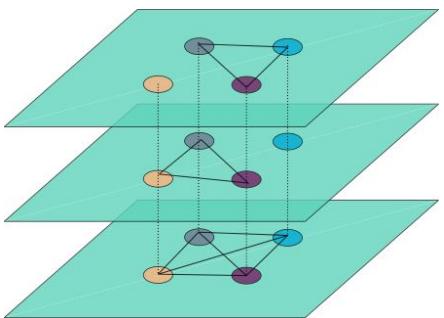
objective/cost  
function



data  
representation  
(reduced)

Recommended reading: "Multi-omics Analysis" chapter by Jonathan Ronen  
from the "Computational Genomics with R"

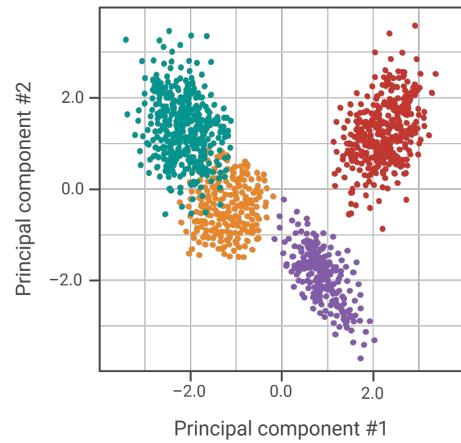
# Primer on dimensionality reduction



data  
representation  
(full)



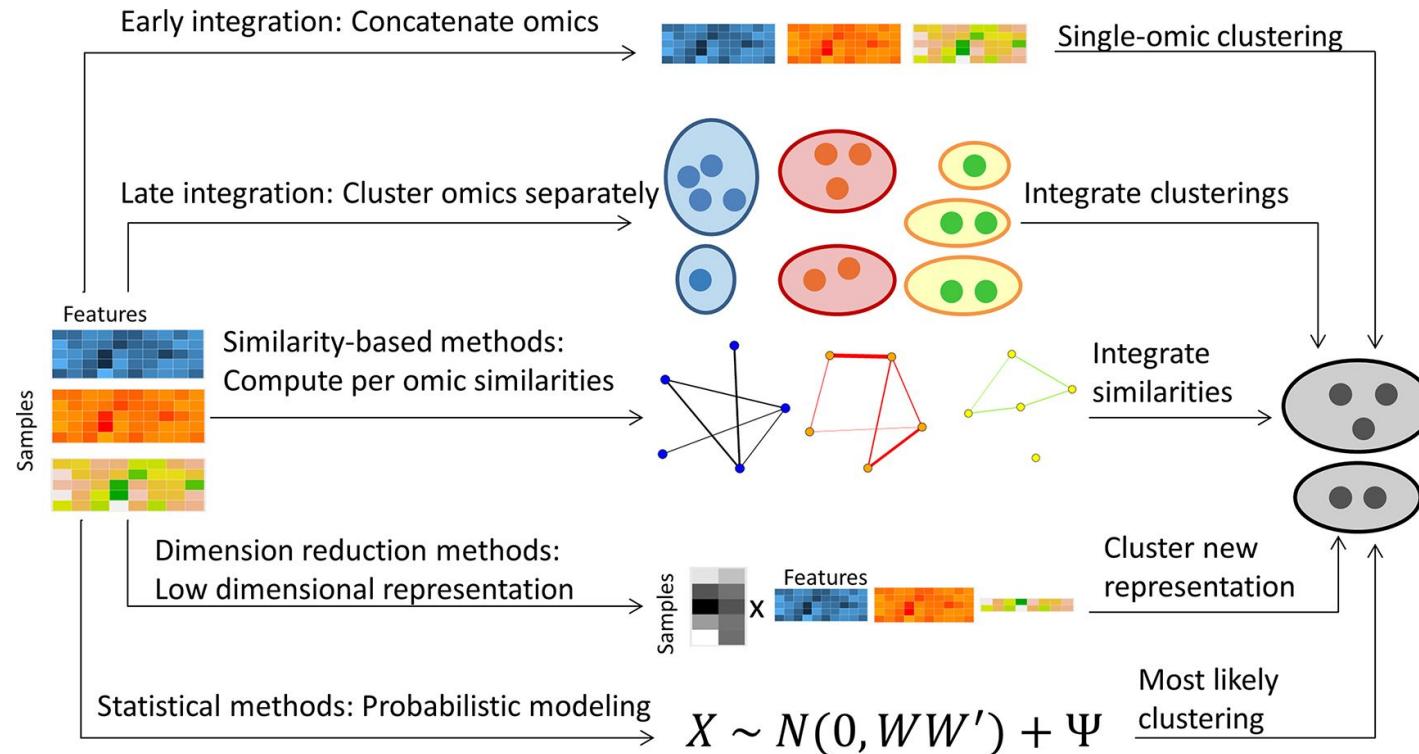
objective/cost  
function



data  
representation  
(reduced)

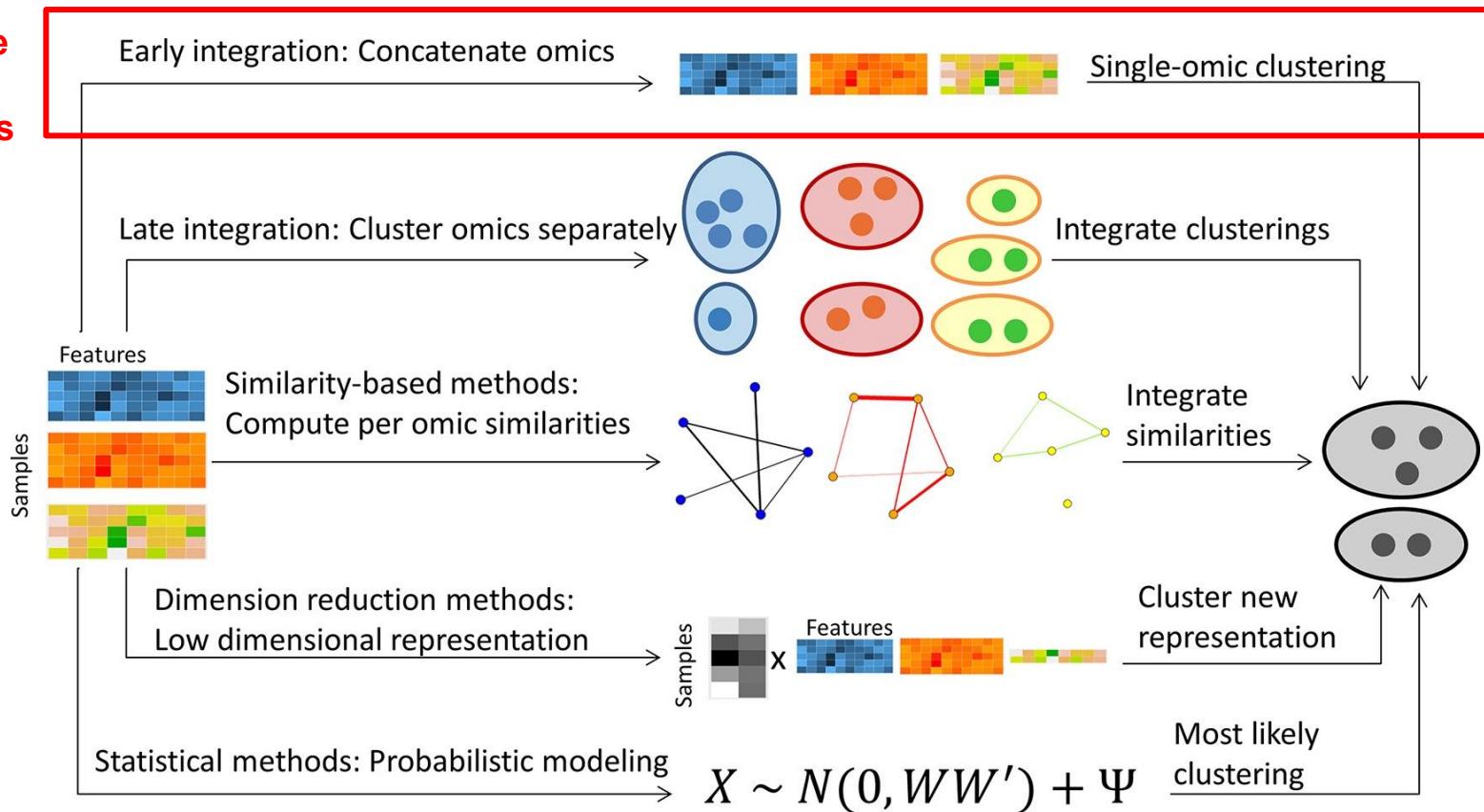
Recommended reading: "Multi-omics Analysis" chapter by Jonathan Ronen  
from the "Computational Genomics with R"

# Approaches to multi-omic data integration

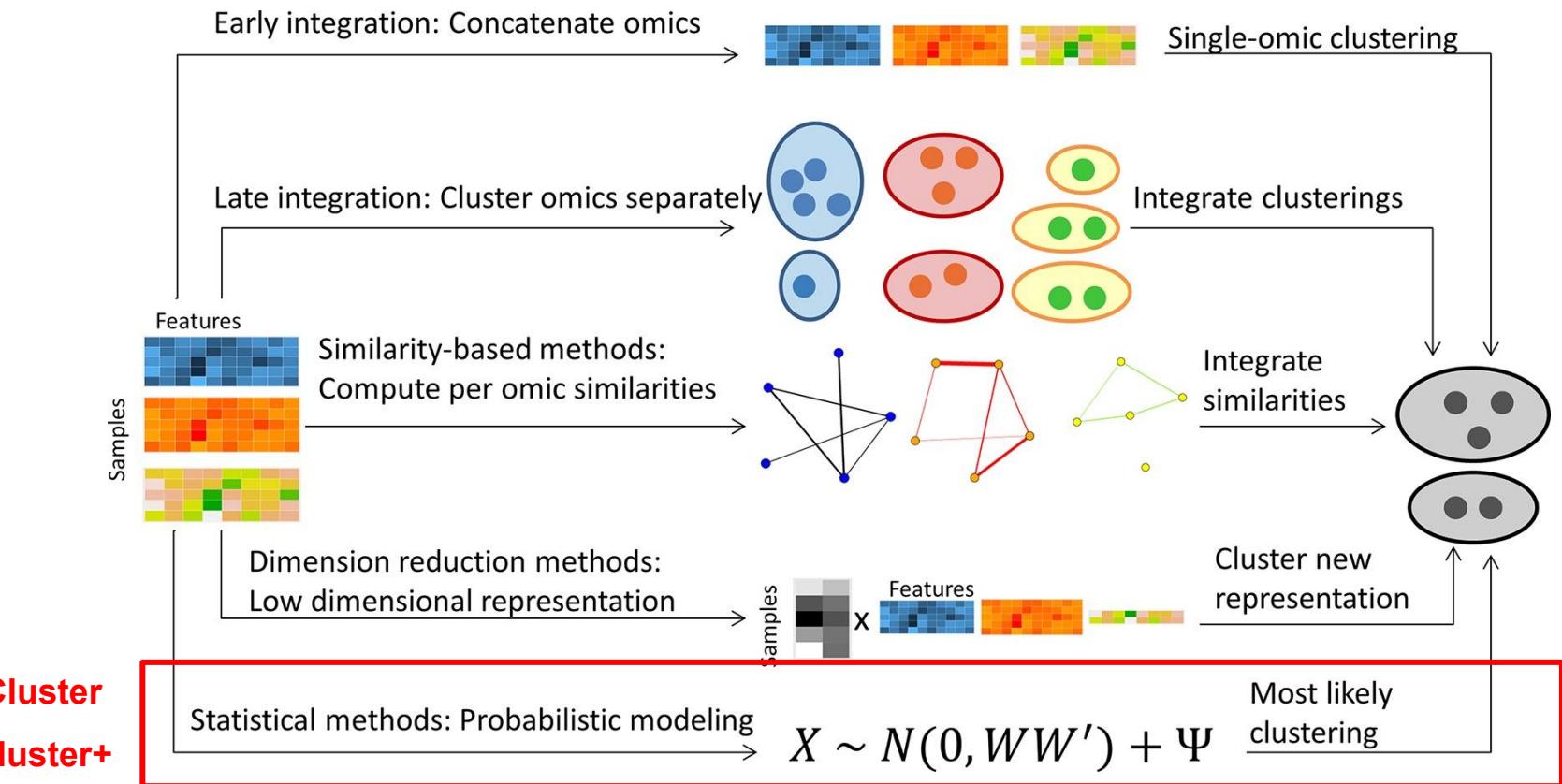


# Methods for multi-omic data integration

## Multiple Factor Analysis (MFA)

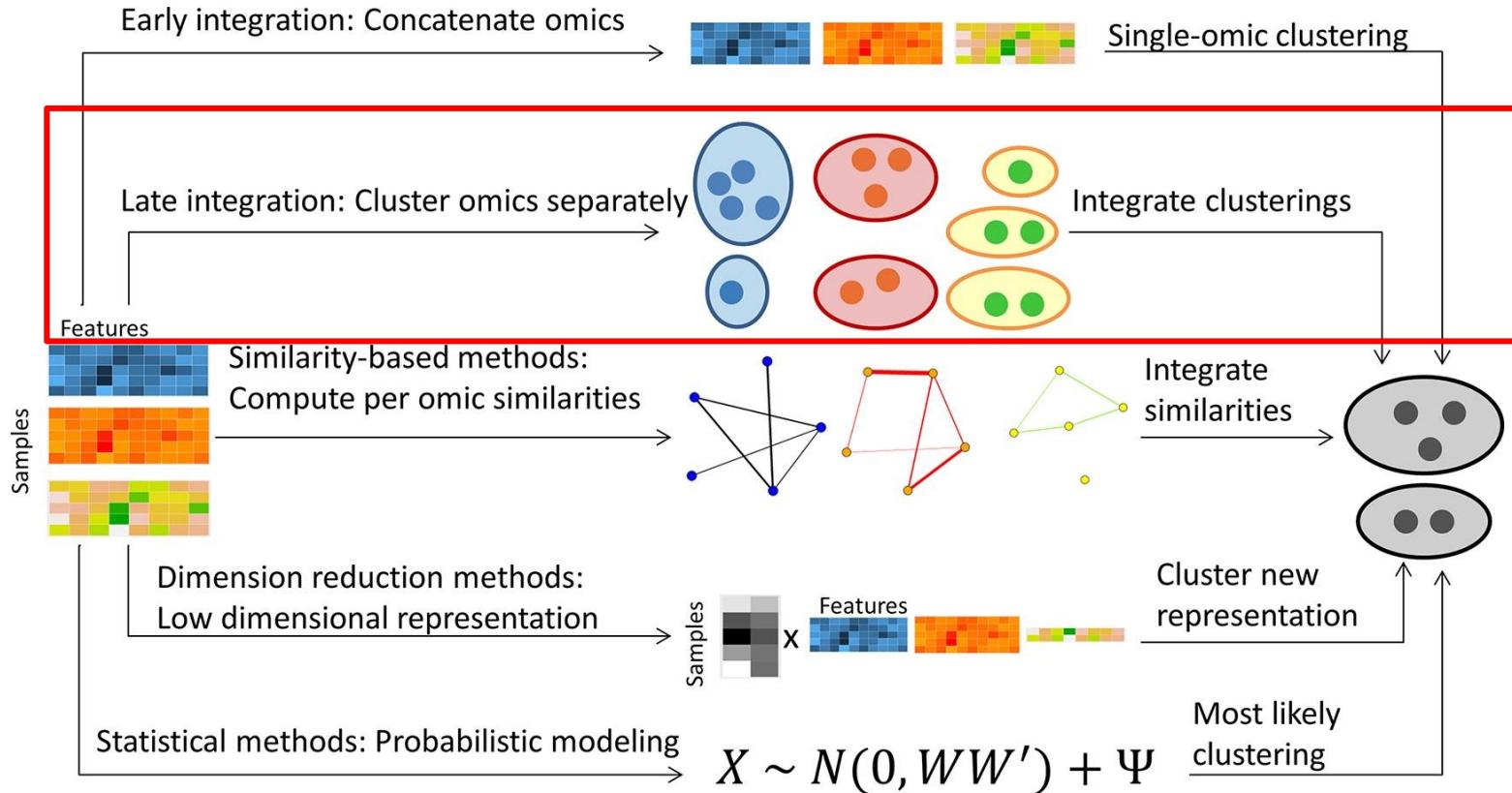


# Methods for multi-omic data integration



# Methods for multi-omic data integration

PINS  
PINS+



# Methods for multi-omic data integration



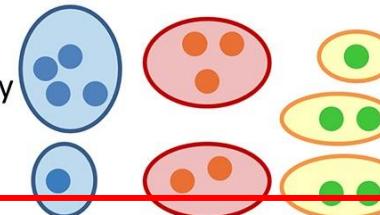
SUMO

Early integration: Concatenate omics



Single-omic clustering

Late integration: Cluster omics separately



Integrate clusterings

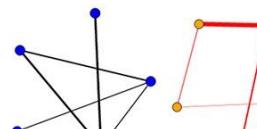
Features



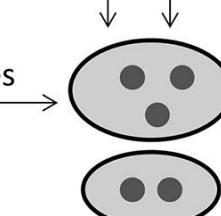
Similarity-based methods:

Compute per omic similarities

Samples



Integrate similarities



Cluster new representation



Features

Dimension reduction methods:  
Low dimensional representation

Statistical methods: Probabilistic modeling

$$X \sim N(0, WW') + \Psi$$

Most likely  
clustering

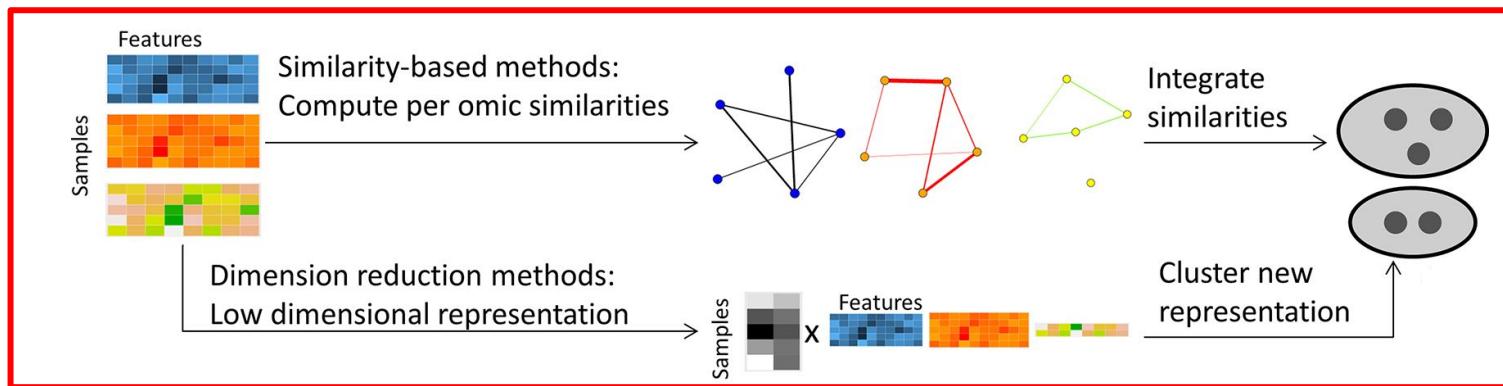
# SUMO: Subtyping Using Multi-Omic data



[github.com/ratan-lab/sumo](https://github.com/ratan-lab/sumo)

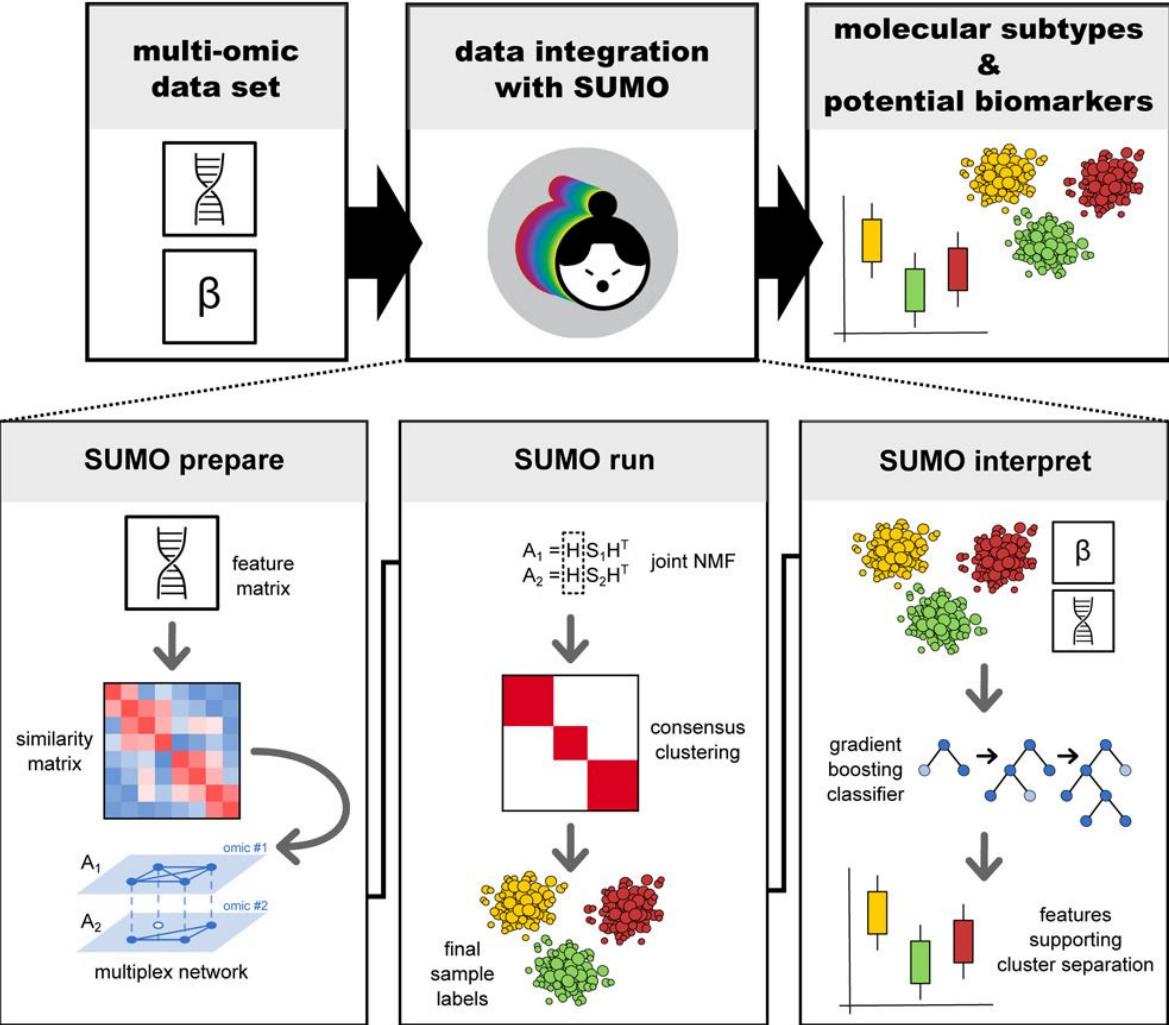


[pypi.org/project/python-sumo](https://pypi.org/project/python-sumo)

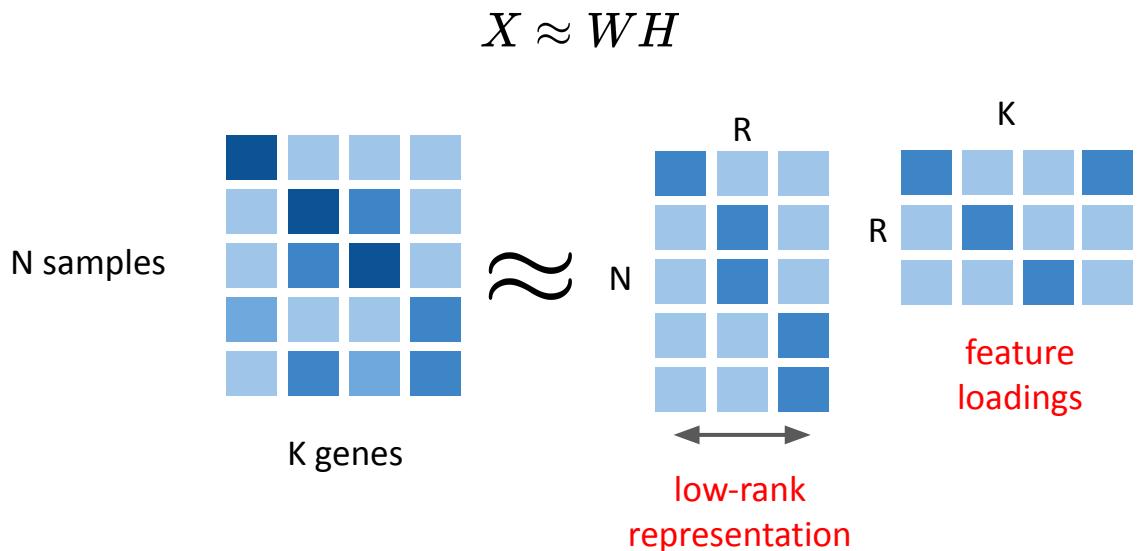


# SUMO workflow

extendable & modular framework  
does not require imputation  
omic-specific similarity metrics:  
- euclidean distance based  
- cosine similarity  
- correlation



# Factorization basics



**NMF**

$$X \approx WH$$

**Symmetric  
NMF**

$$A \approx HH^T$$

# SUMO: Factorization

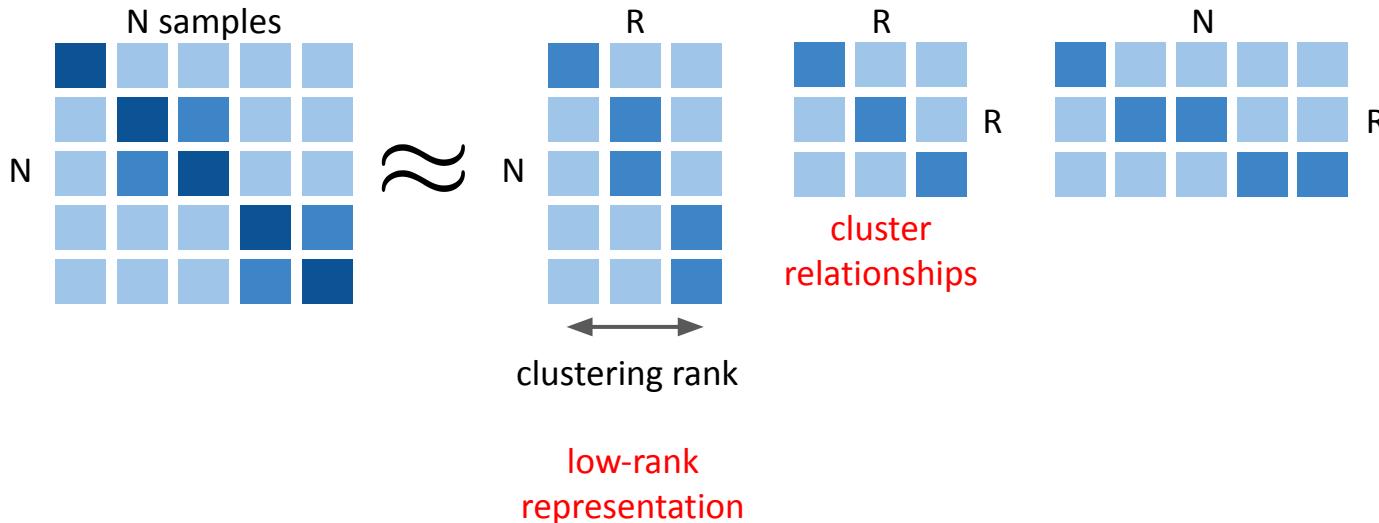
NMF

$$X \approx WH$$

Symmetric  
NMF

$$A \approx HH^T$$

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



# SUMO: Factorization

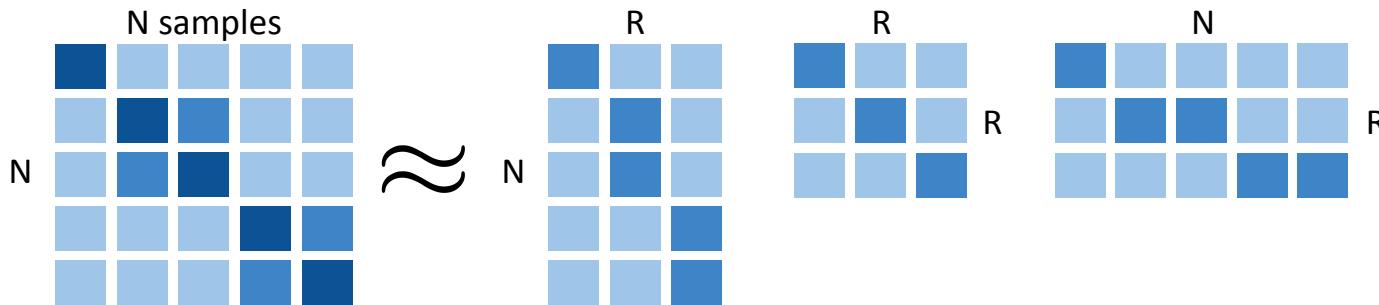
$$A^i \approx H S_i H^T$$

**NMF**

$$X \approx WH$$

**Symmetric  
NMF**

$$A \approx HH^T$$

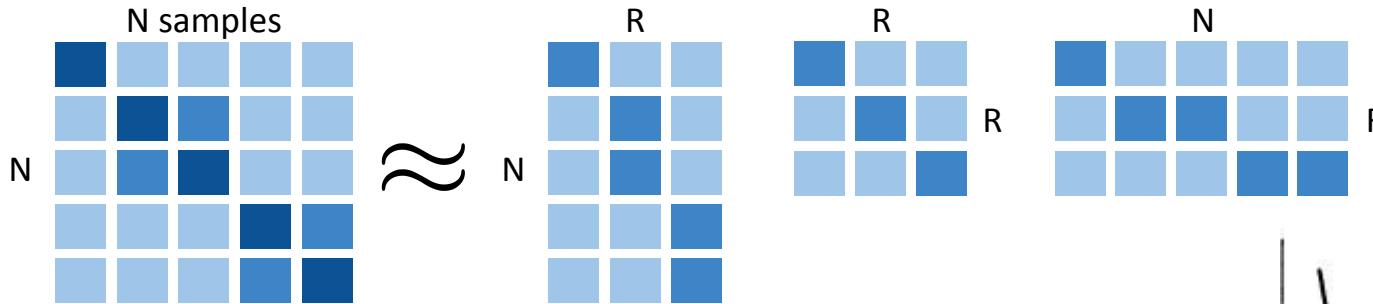


$$\min_{H, S_i \geq 0} \sum_i \lambda_i \left\| W_i \circ (A_i - HS_i H^T) \right\|_F^2 + \eta \|H\|_F^2$$

different number of samples      missing data      sparsity

# SUMO: Factorization

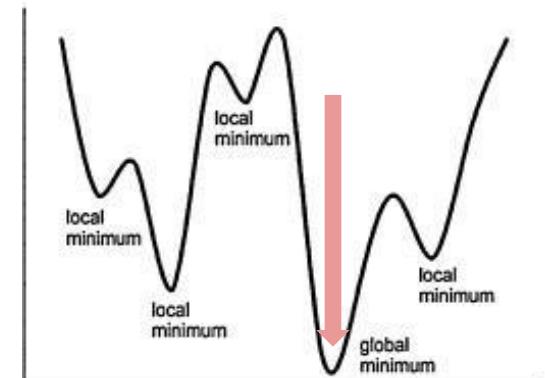
$$A^i \approx H S_i H^T$$



$$\min_{H, S_i \geq 0} \sum_i \lambda_i \left\| W_i \circ (A_i - HS_i H^T) \right\|_F^2 + \eta \|H\|_F^2$$

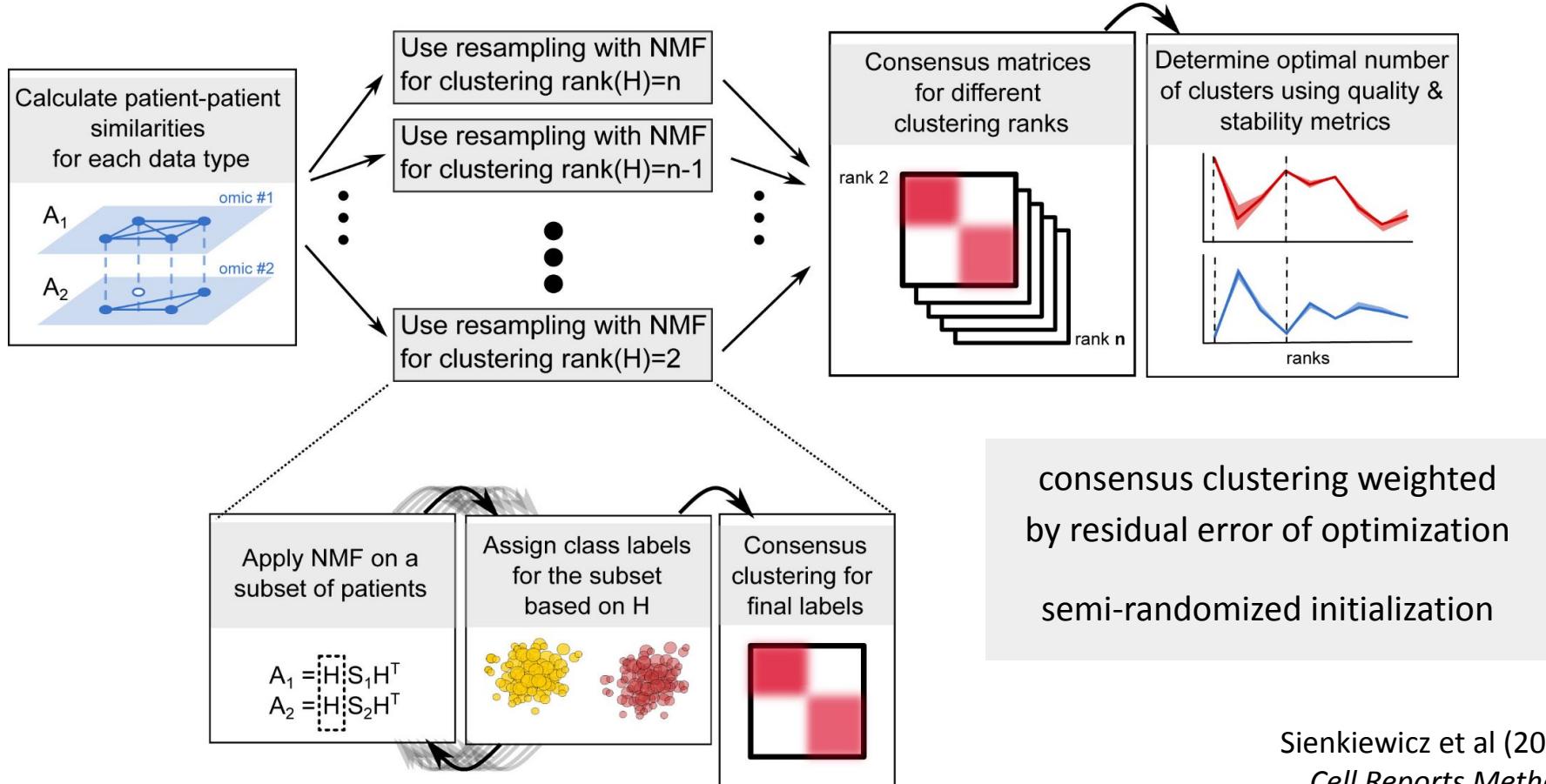
different number of samples      missing data      sparsity

<b>NMF</b>	$X \approx WH$
<b>Symmetric NMF</b>	$A \approx HH^T$

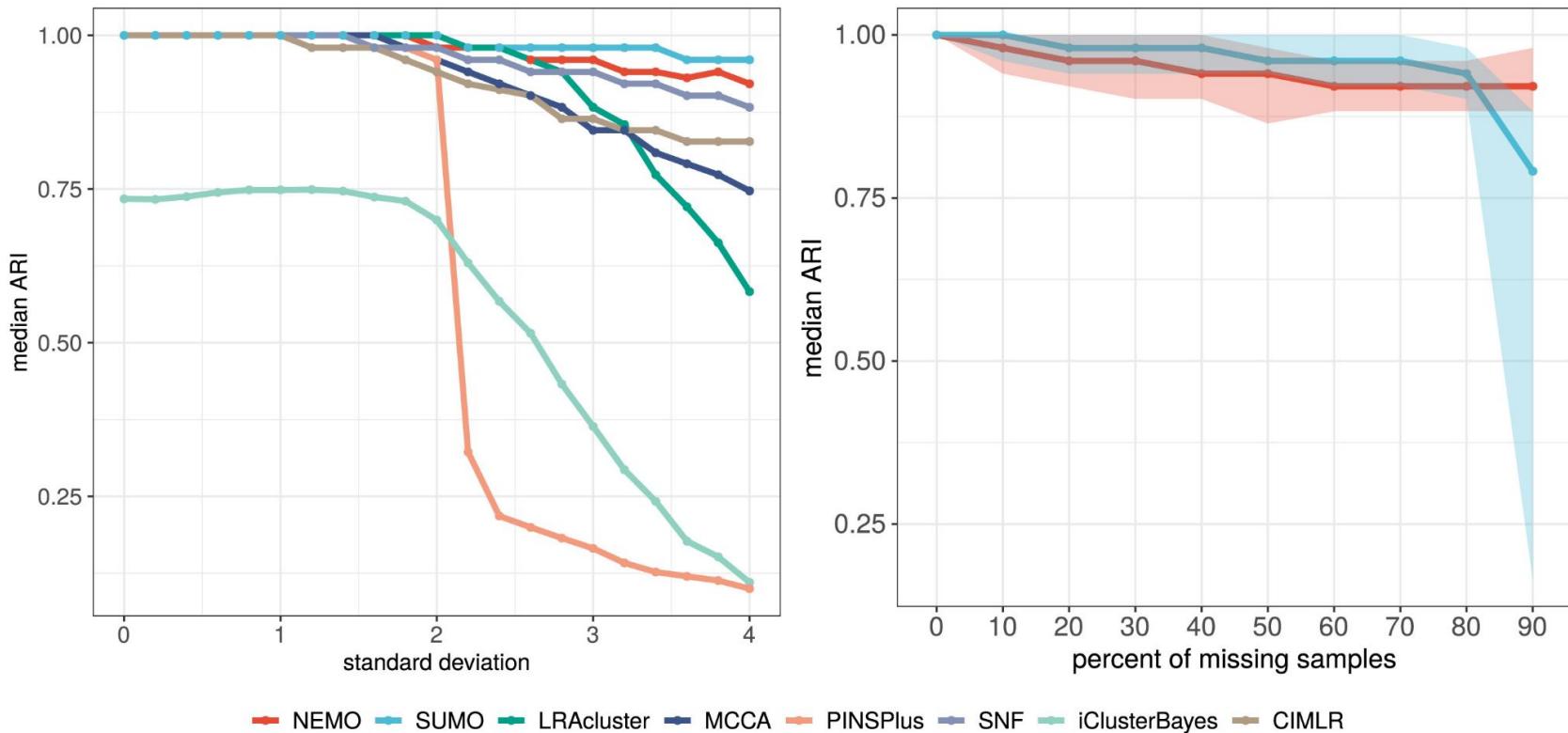


using gradient descent

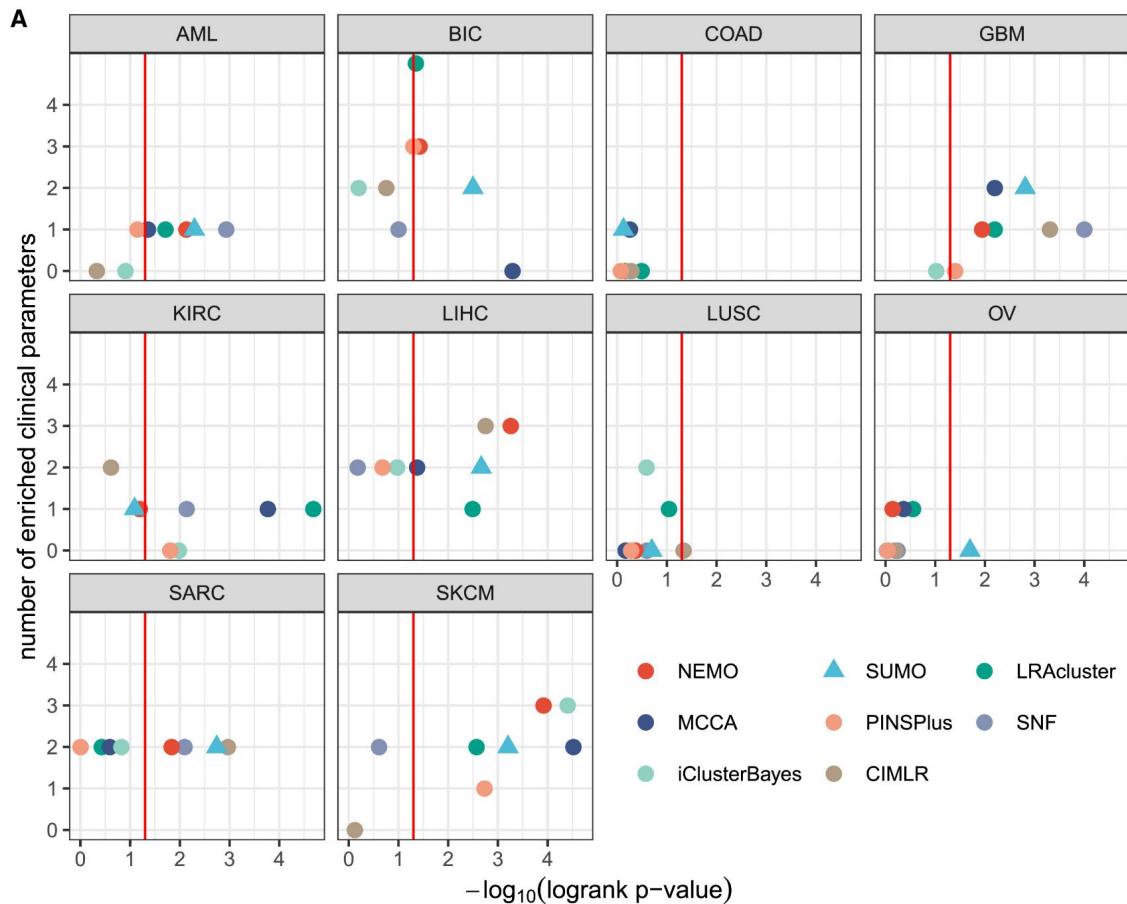
# SUMO: Re-sampling based approach



# SUMO improves performance with noisy and incomplete data

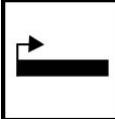
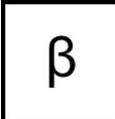
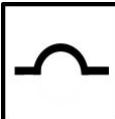


SUMO consistently identifies differential subtypes in TCGA data using an established benchmark

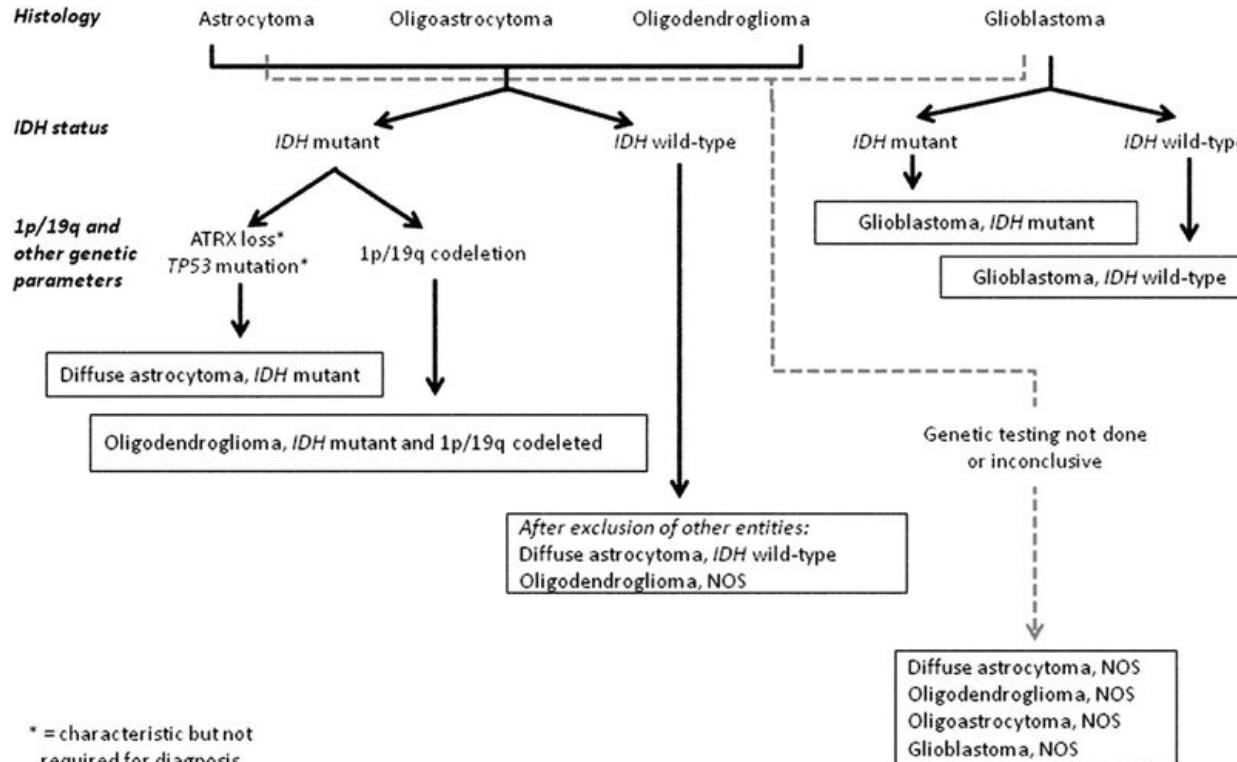


# SUMO detects latent relationships between patients in Lower-Grade Gliomas

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

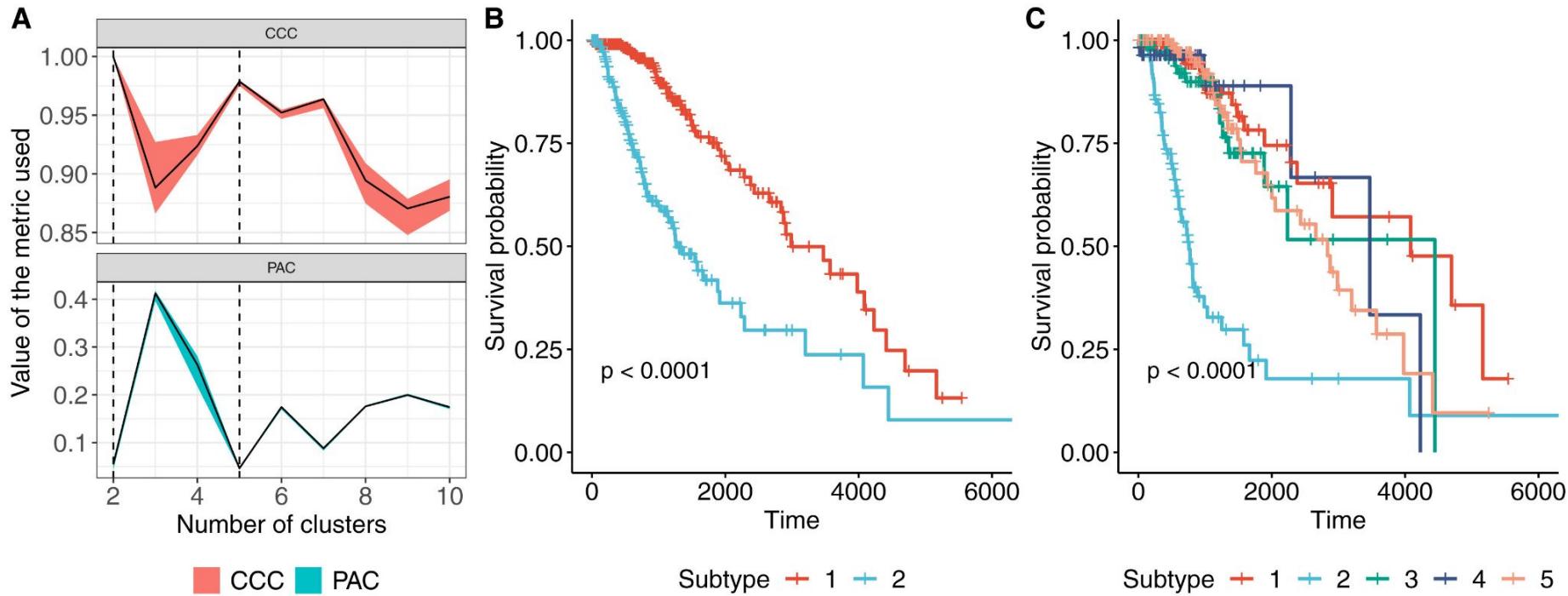
omic	# samples	tech
 RNA-seq	530	Illumina HiSeq
 DNA methylation	530	Illumina Infinium HumanMethylation450
 miRNA expression	524	Illumina HiSeq

# WHO guidelines for Low-Grade Gliomas

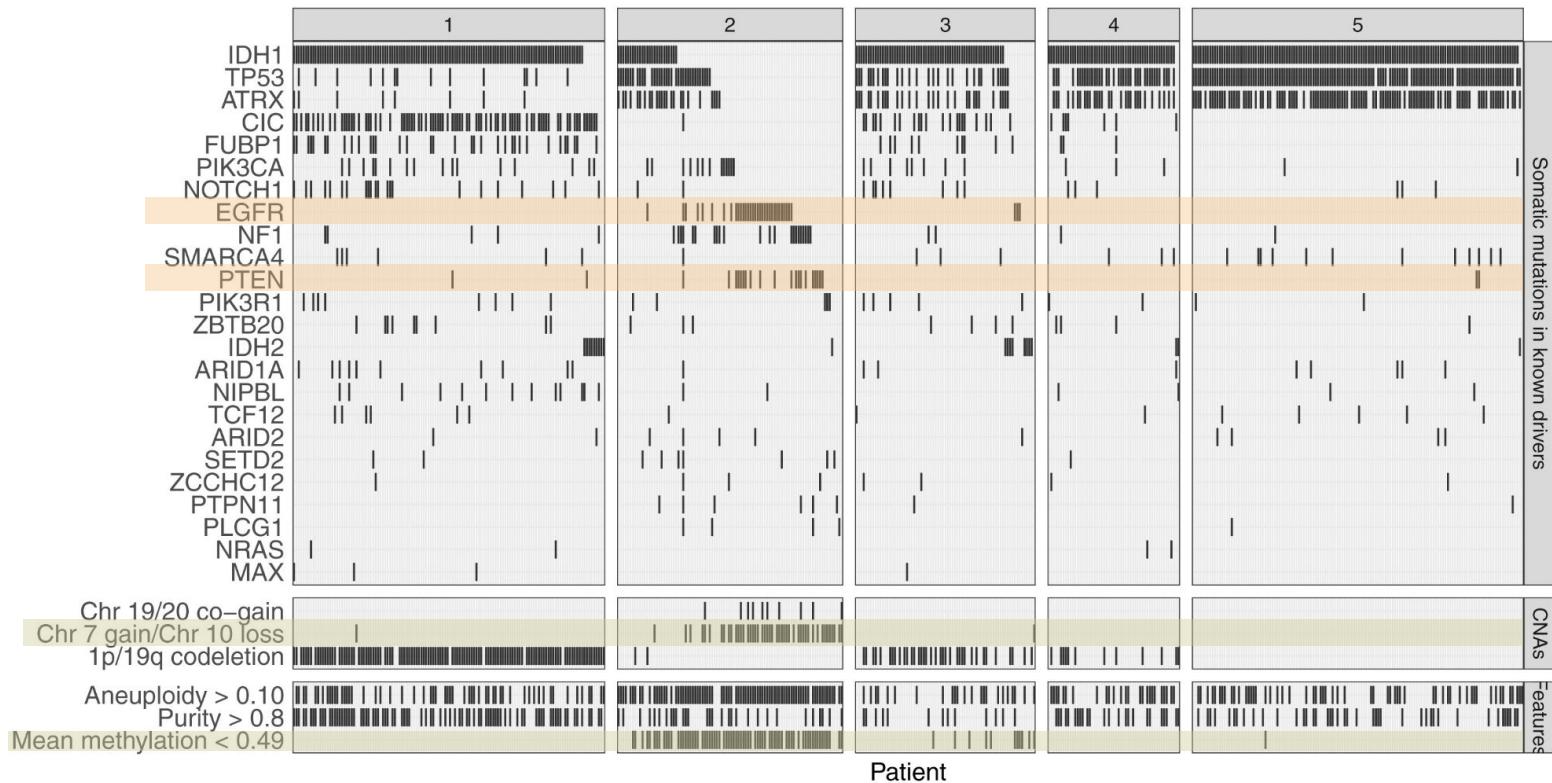


- Gliomas are neuroepithelial tumors originating from glial cells
- Postoperative RT remains the standard of care for adult grade 2/3
- ~80% have mutations in IDH1, 90% of those are R132H.
- ~50% have mutations in TP53, 40% have mutations in ATRX and 20% in CIC

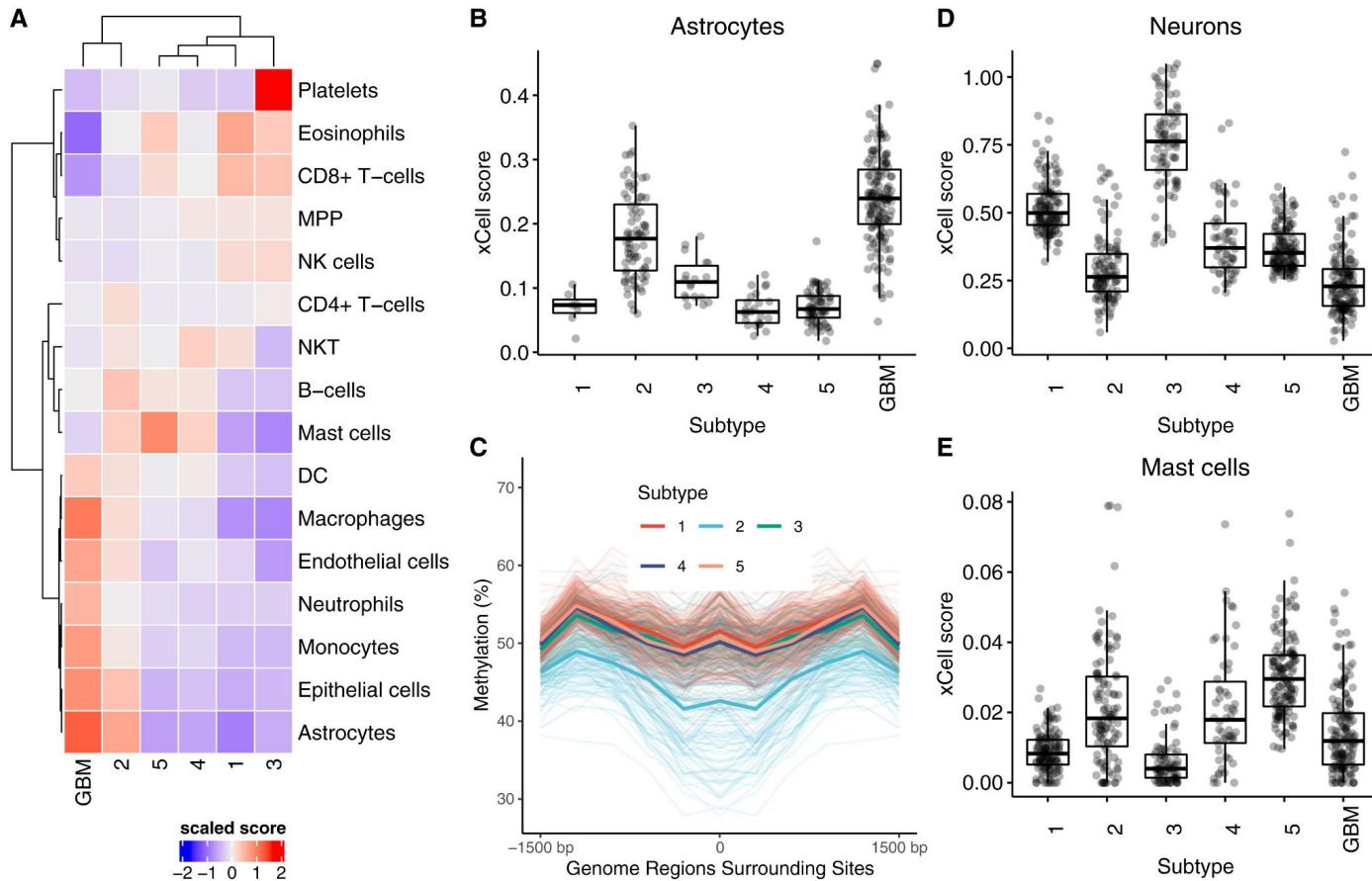
# SUMO identifies a subtype of LGG with differential prognosis & GBM-like features



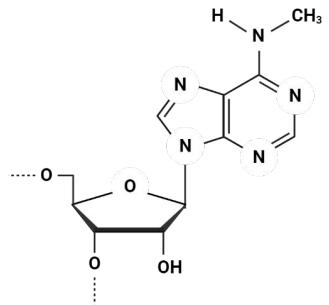
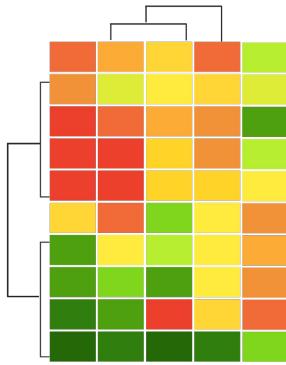
# SUMO identifies a subtype of LGG with differential prognosis & GBM-like features



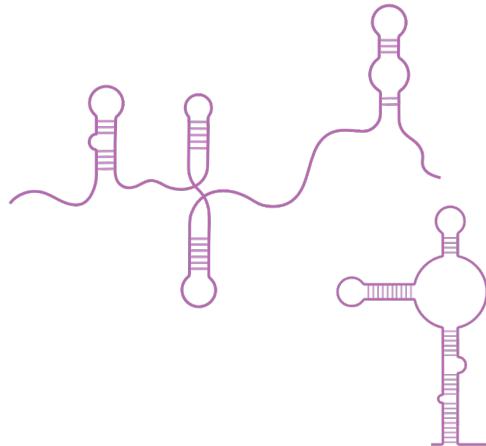
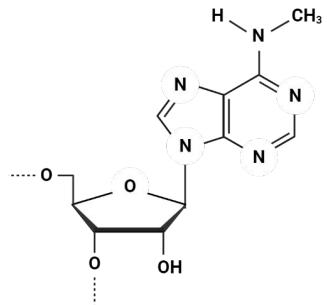
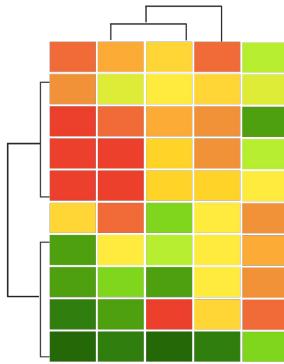
SUMO  
subtype 2  
cellular profile  
is more  
similar to  
GBMs than  
other LGGs



# Single-molecule multi-omics?

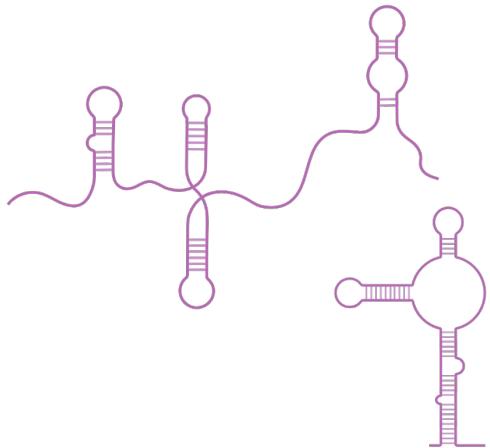
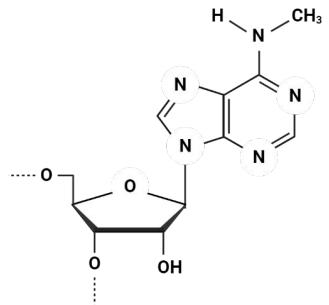
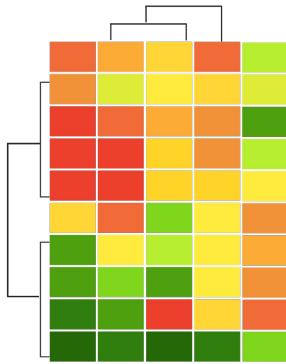


# Single-molecule multi-omics?



- matching transcript expression with base modifications
- multiple distinct base modification

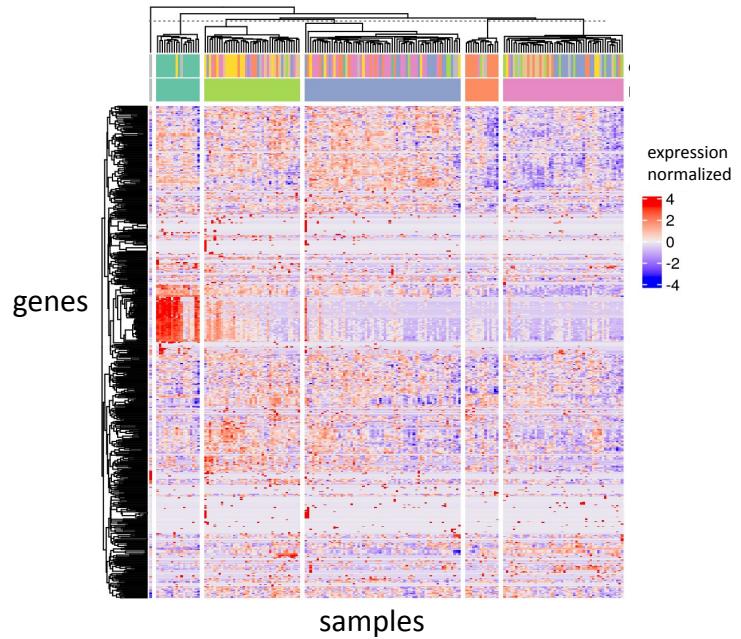
# Single-molecule multi-omics?



- matching transcript expression with base modifications
- multiple distinct base modification
- transcripts as features vs transcripts as samples
- secondary structure prediction and validation

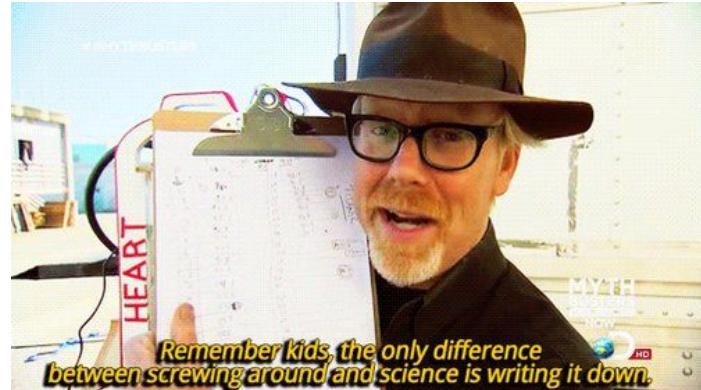
# Conclusions

- Heatmaps are your friend



# Conclusions

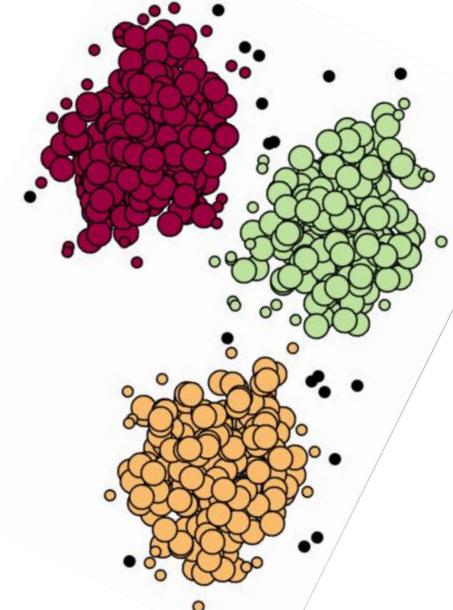
- Heatmaps are your friend
- Analysis-aware experimental design is essential
- Keep track of the metadata



Adam Savage

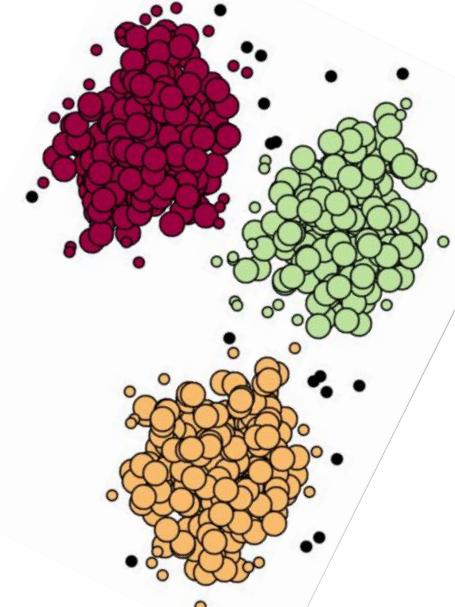
# Conclusions

- 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
- Single-molecule omics are on the way!



## Acknowledgements

### SUMO Team

**University of Virginia**

Aakrosh Ratan

Ajay Chatrath

Nathan Sheffield

John Lawson

**National University of Singapore**

Jinyu Chen

Louxin Zhang

**WCM**

Mason Lab

Melnick Lab

### SUMO Funding



**UNIVERSITY  
of VIRGINIA**

**CENTER for  
PUBLIC  
HEALTH  
GENOMICS**



National Institutes  
of Health



NIGMS

**NATIONAL  
RESEARCH  
FOUNDATION**  
SINGAPORE

**Tri-Institutional PhD Program  
Computational Biology & Medicine**



Memorial Sloan Kettering  
Cancer Center



**Weill Cornell  
Medicine**