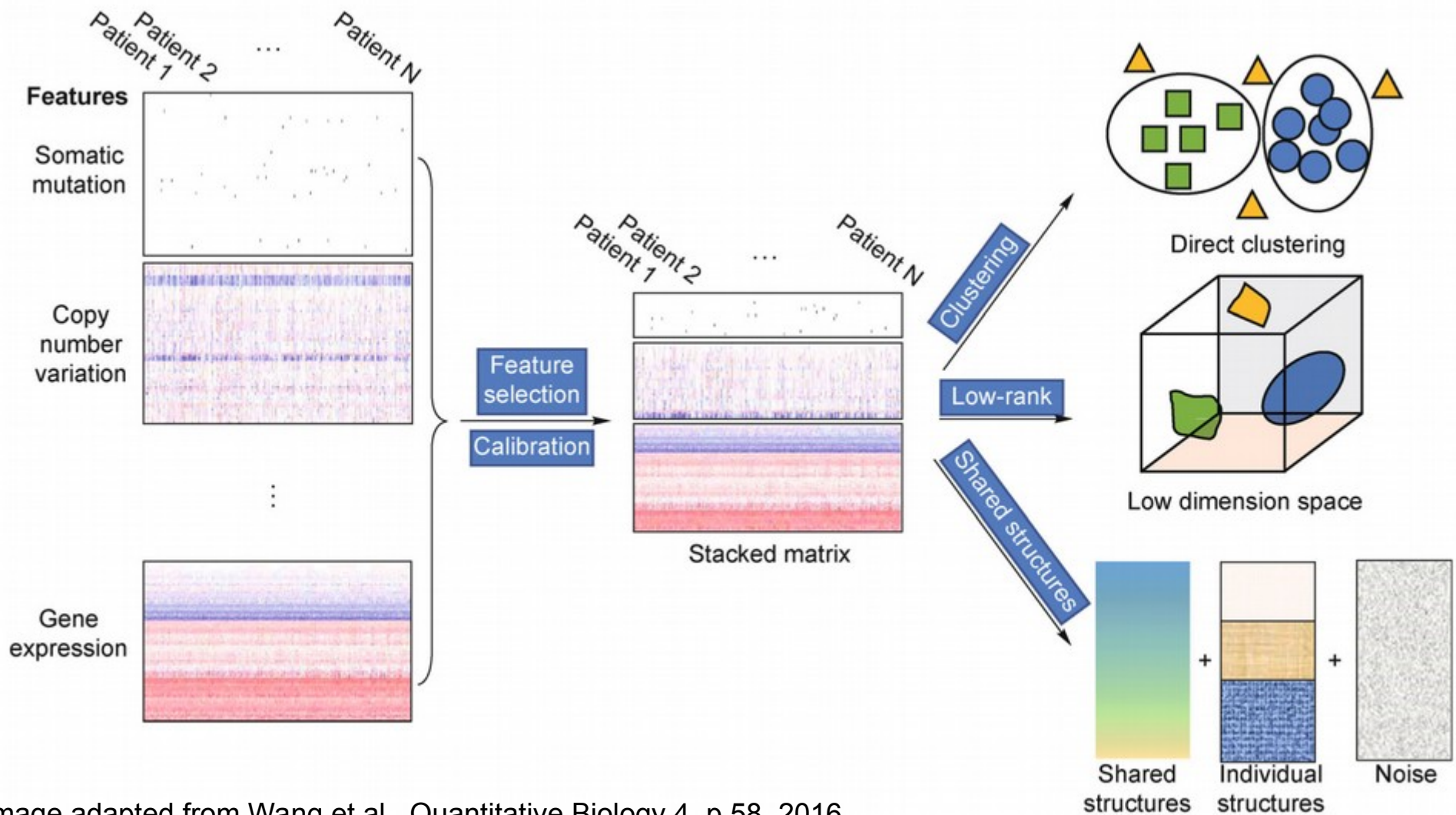# Unsupervised OMICs Integration
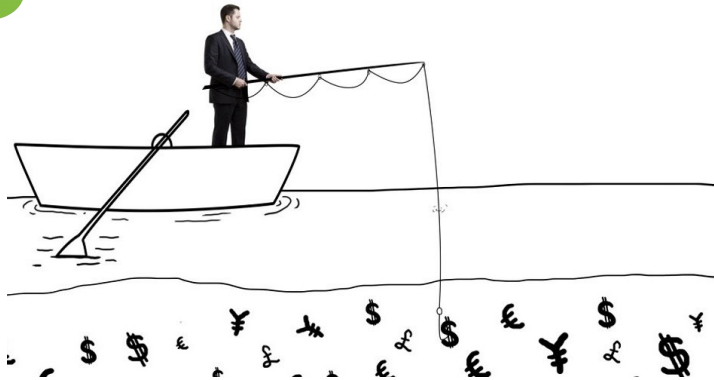
ISMB / ECCB 2021
Tutorial: a practical introduction to multi-omics integration and network analysis
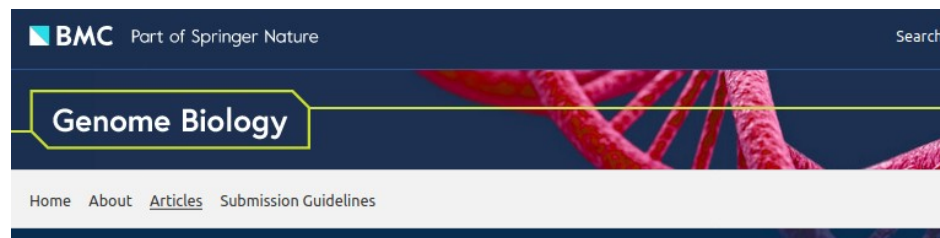Nikolay Oskolkov, NBIS SciLifeLab, 22.07.2021

Image adapted from Wang et al., Quantitative Biology 4, p.58, 2016

# Find Something in My Data

**Fishing expedition**

- I do not understand your biological hypothesis

- I do not have any



## BMC Part of Springer Nature

### Genome Biology

Home   About   Articles   Submission Guidelines

Editorial | Open Access | Published: 03 September 2020

## A hypothesis is a liability

Itai Yanai ✉ & Martin Lercher ✉

*Genome Biology* **21**, Article number: 231 (2020) | Cite this article
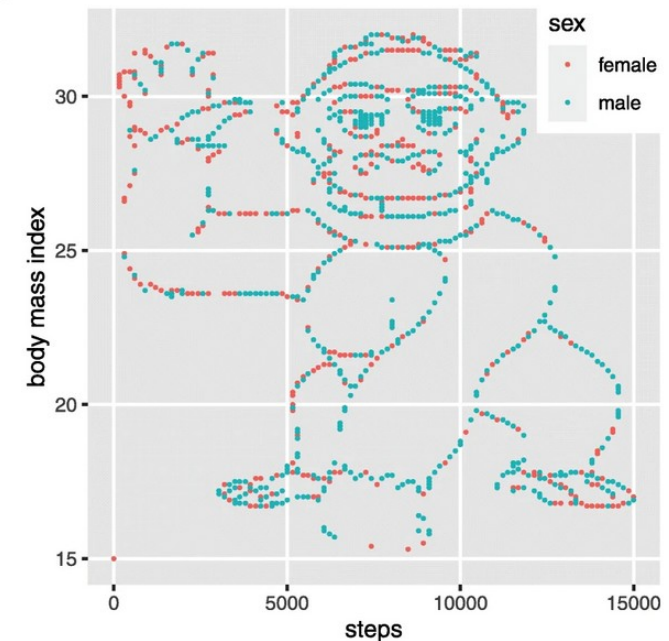
12k Accesses | 619 Altmetric | Metrics

" 'When someone seeks,' said Siddhartha, 'then it easily happens that his eyes see only the thing that he seeks, and he is able to find nothing, to take in nothing. [...] Seeking means: having a goal. But finding means: being free, being open, having no goal.' " Hermann Hesse

There is a hidden cost to having a hypothesis. It arises from the relationship between night science and day science, the two very distinct modes of activity in which scientific ideas are generated and tested, respectively [1, 2]. With a hypothesis in hand, the impressive strengths of day science are unleashed, guiding us in designing tests, estimating parameters, and throwing out the hypothesis if it fails the tests. But when we analyze the results of an experiment, our mental focus on a specific hypothesis can prevent us from exploring other aspects of the data, effectively blinding us to new ideas. A hypothesis then becomes a liability for any night science explorations. The corresponding limitations on our creativity, self-imposed in hypothesis-driven research, are of particular concern in the context of modern biological datasets, which are often vast and likely to contain hints at multiple distinct and potentially exciting discoveries. Night science has its own liability though, generating many spurious relationships and false hypotheses. Fortunately, these are exposed by the light of day science, emphasizing the complementarity of the two modes, where each overcomes the



a An artificial dataset given to students with and without explicit hypotheses on the relationship between BMI and the steps taken on a particular day, for men and women. b A plot of the dataset. c The contingency table for students in the two groups ("hypothesis-focused," "hypothesis-free") that discovered the gorilla or not [6]
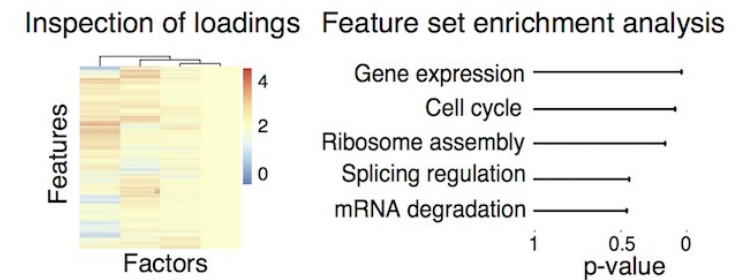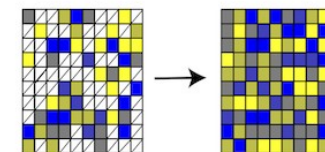
# MOFA: Hypothesis-Free Exploratory Tool
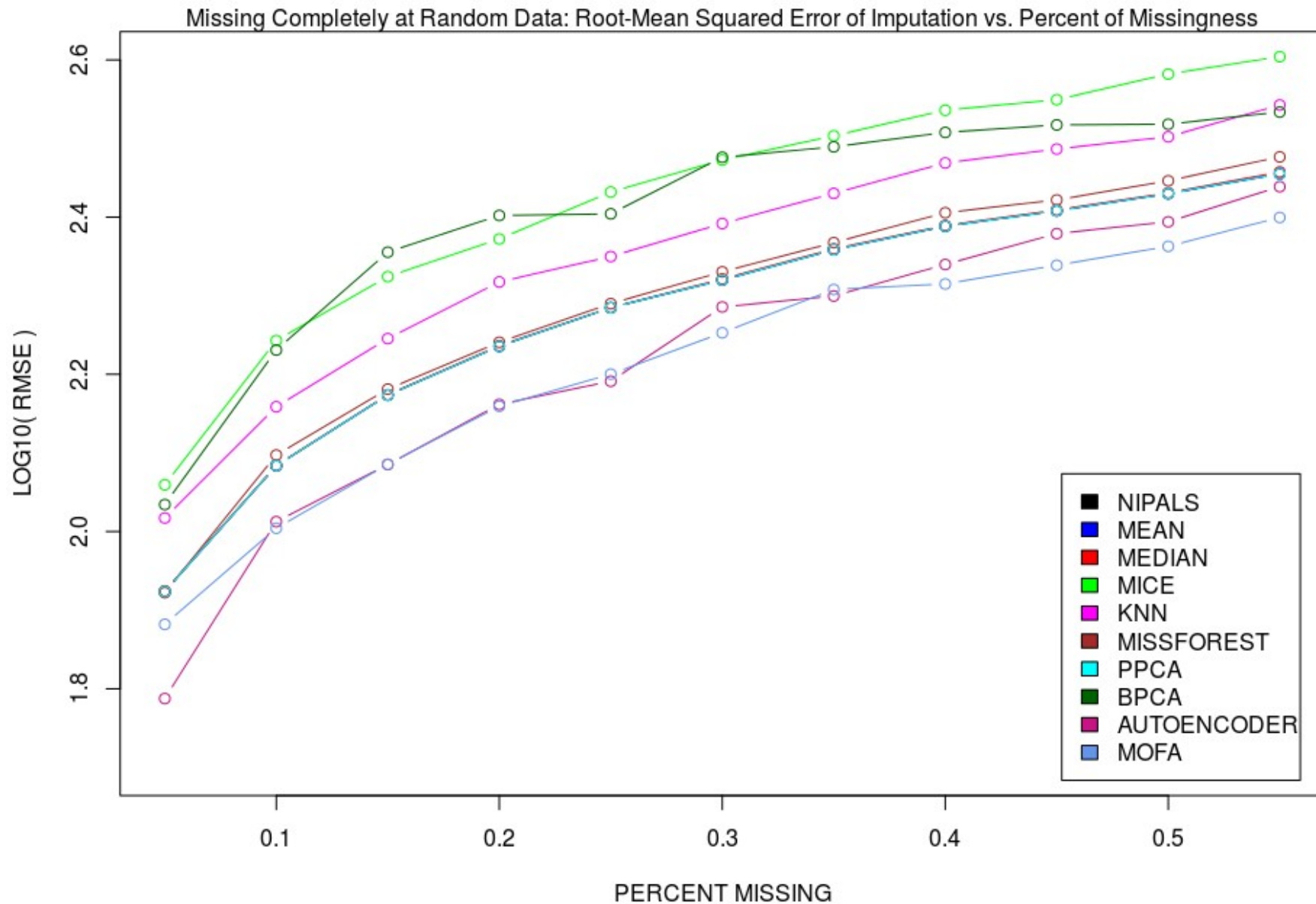


- Visualisation of samples in factor space
- Annotation of factors using (gene set) enrichment analysis
- Imputation of missing values
- Support of OMICs with non-Gaussian distribution including binary and count data

# MOFA: Imputation of Missing Values



Bayesian framework is insensitive to missing data, priors compensate for the lack of data

MOFA: Imputation of Missing Values

# scNMT Data Set:
## scRNAseq + scBSseq + scATACseq

**ARTICLE**

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

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

Parallel single-cell sequencing protocols represent powerful methods for investigating regulatory relationships, including epigenome-transcriptome interactions. Here, we report a single-cell method for parallel chromatin accessibility, DNA methylation and transcriptome profiling. scNMT-seq (single-cell nucleosome, methylation and transcription sequencing) uses a GpC methyltransferase to label open chromatin followed by bisulfite and RNA sequencing. We validate scNMT-seq by applying it to differentiating mouse embryonic stem cells, finding links between all three molecular layers and revealing dynamic coupling between epigenomic layers during differentiation.
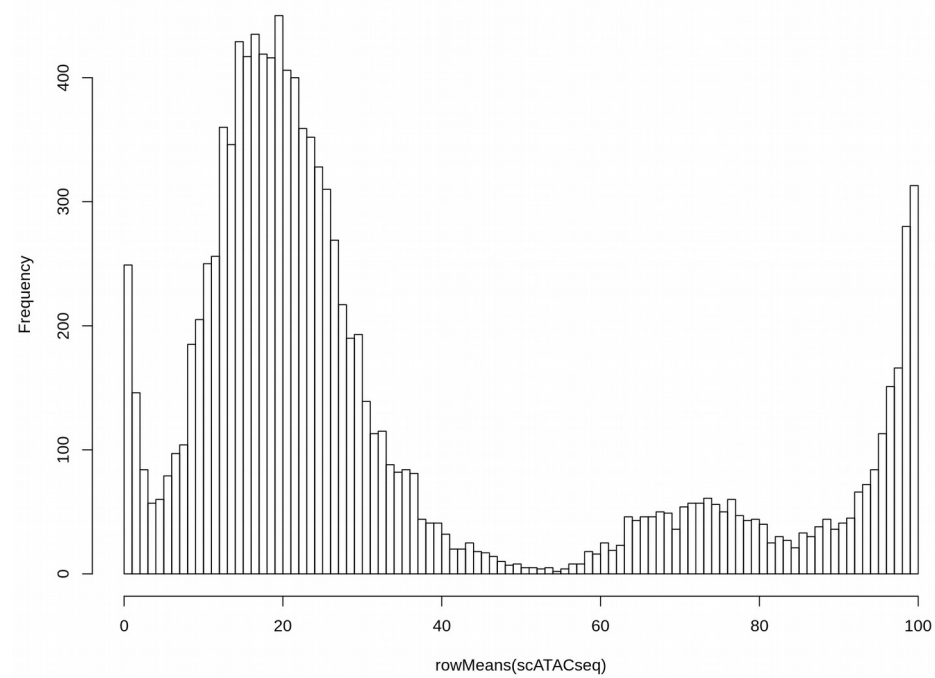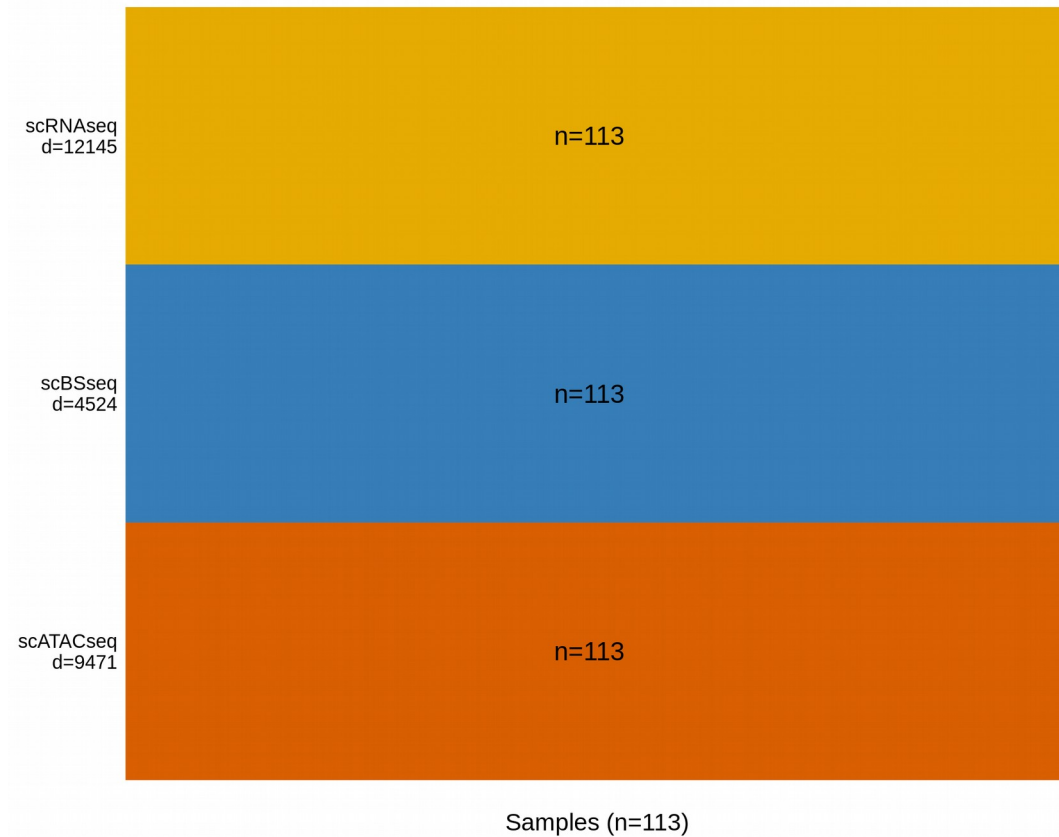
**a**

Lysis / GpC methylase

GpC methylase treatment

SMART-Seq2 — Sequencing and mapping

scBS — Sequencing, mapping and splitting

Transcriptome | DNA methylation | Accessibility
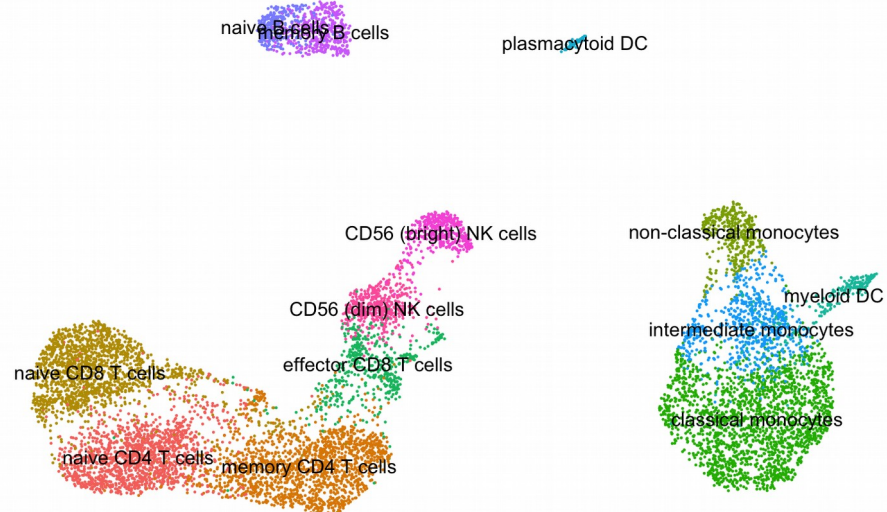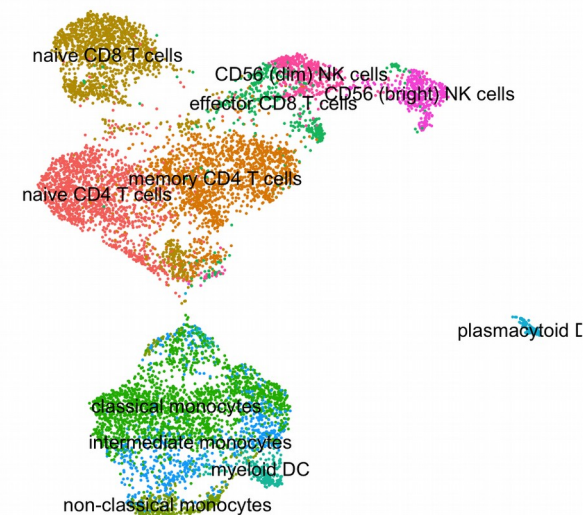
scNMT Data Set: Distributions

scNMT Data Set: Summary Stats

ESC and EB cells are separable on the heatmap built on loadings of the MOFA latent factors

MOFA for scOmics Integration: 10X PBMC

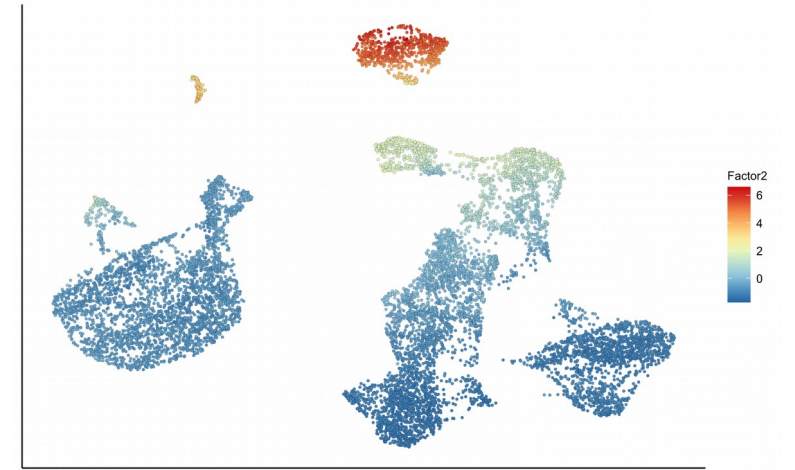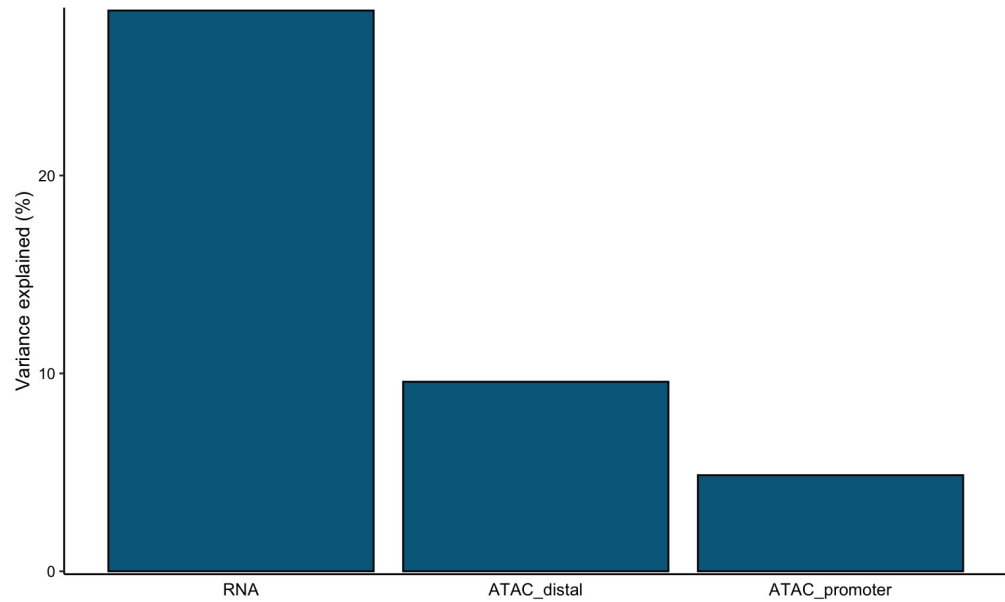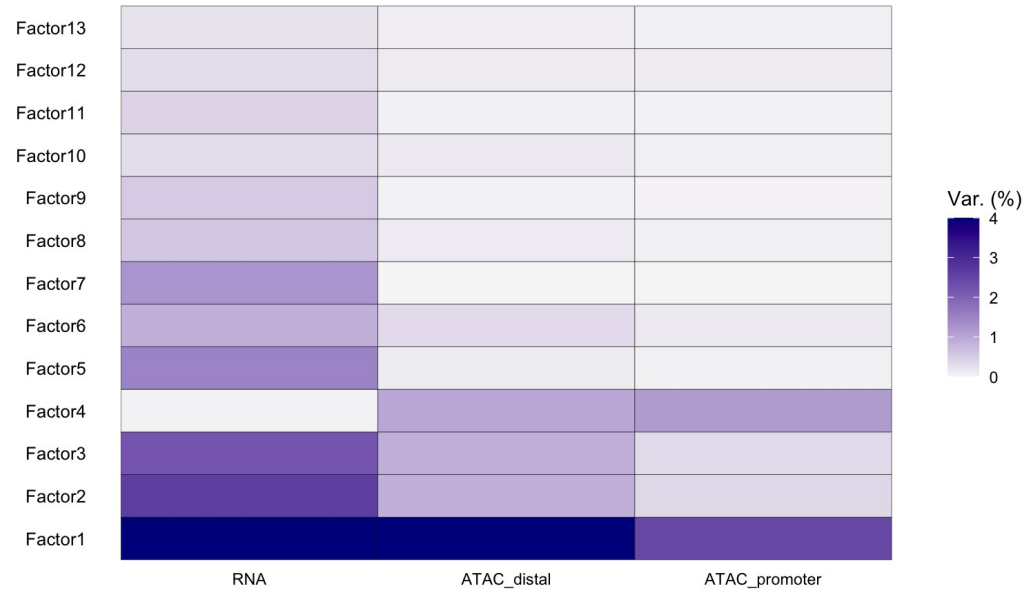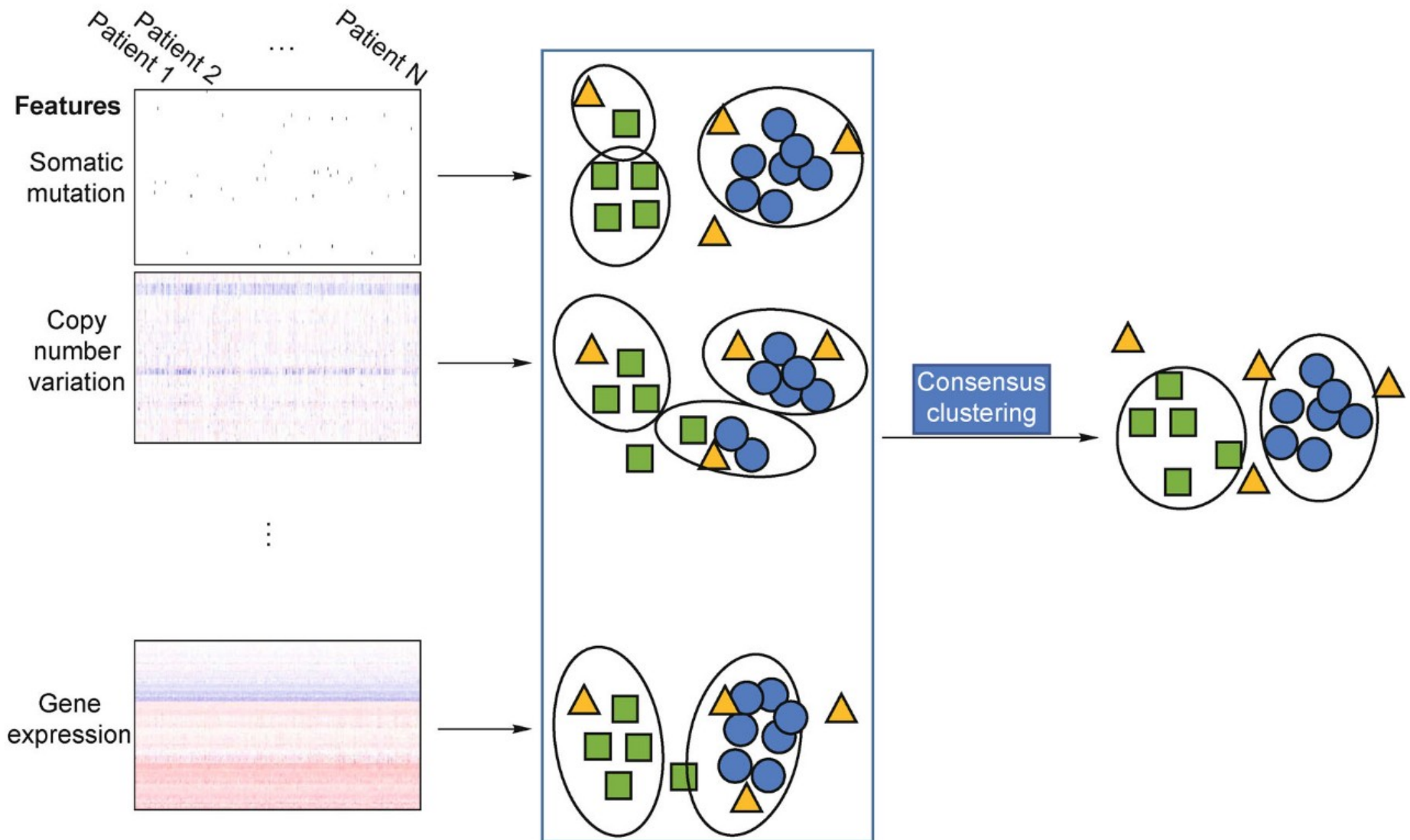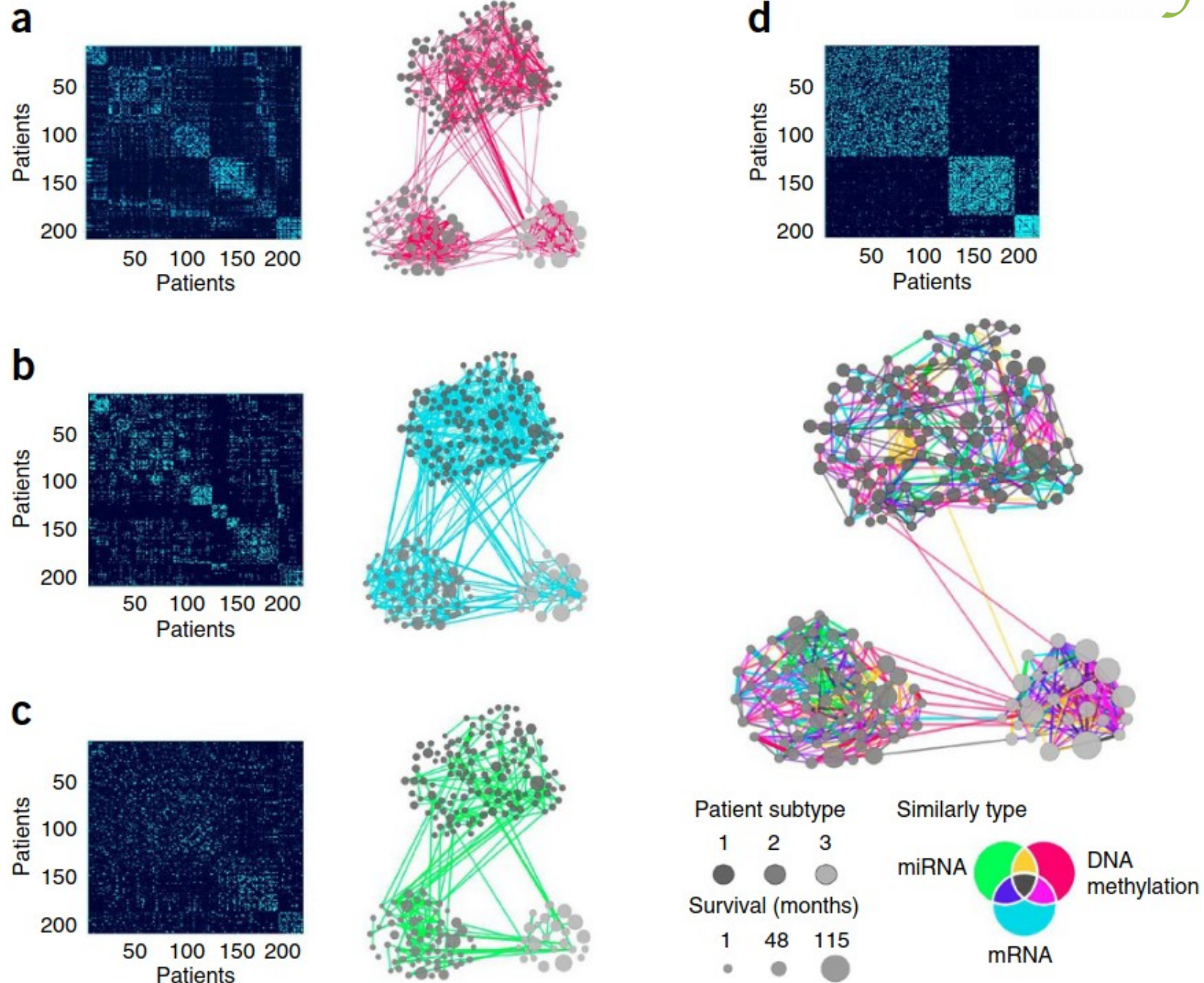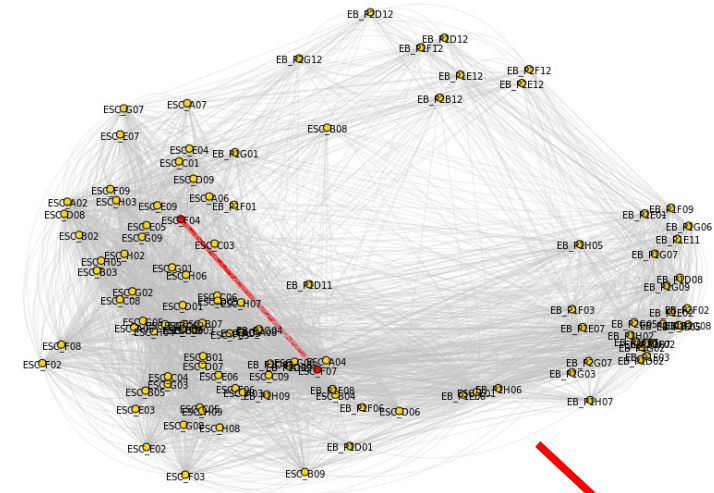Other Unsupervised
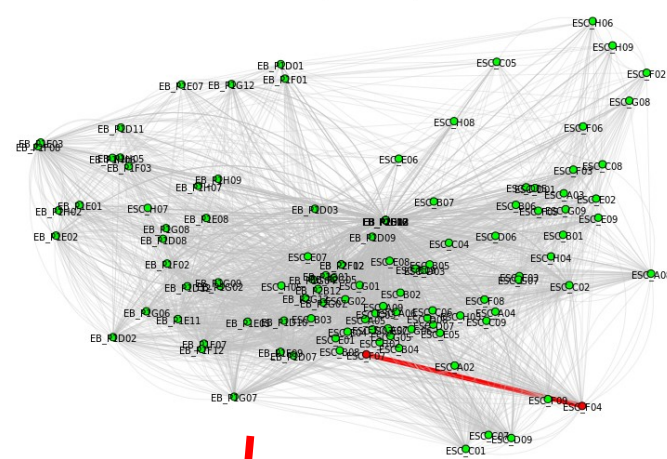Integrative OMICs Methods

**Figure 2. Clustering of clusters.** This kind of methods first clusters in every single omics dataset and then integrates the primary clustering results into final cluster assignments.
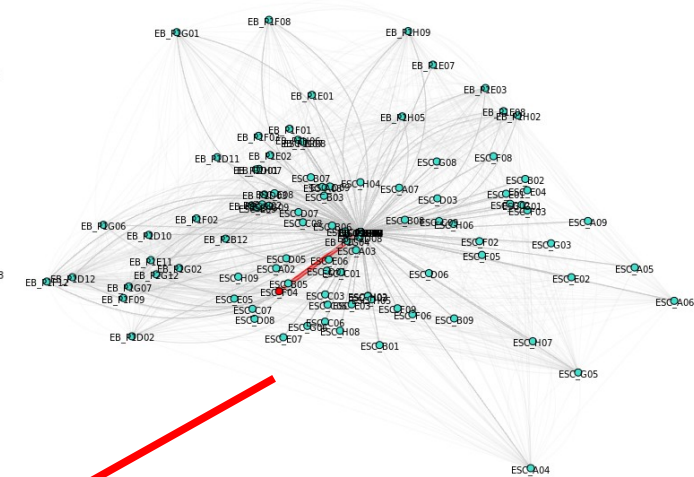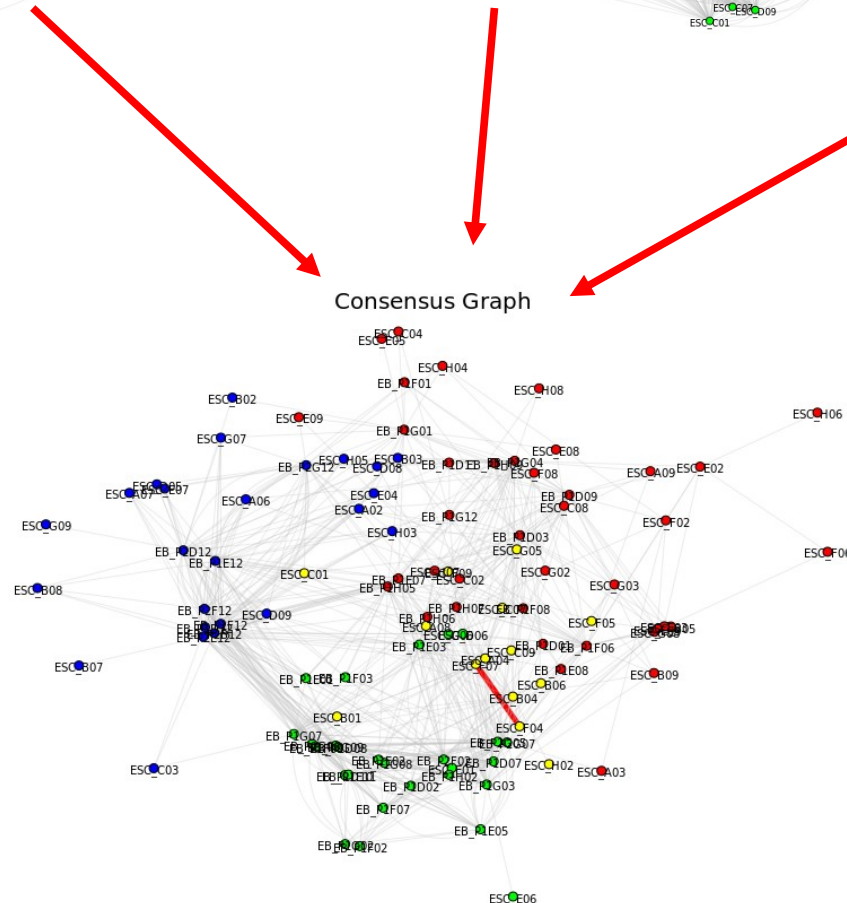
# Similarity Network Fusion (SNF)



Patient subtype
1   2   3

Survival (months)
1   48   115

Similarly type

miRNA        DNA methylation

mRNA

Wang, B., Mezlini, A., Demir, F. et al. . Nat Methods 11, 333–337 (2014)

# Graph-Based OMICs Integration



Keep edges consistently present across the OMICs
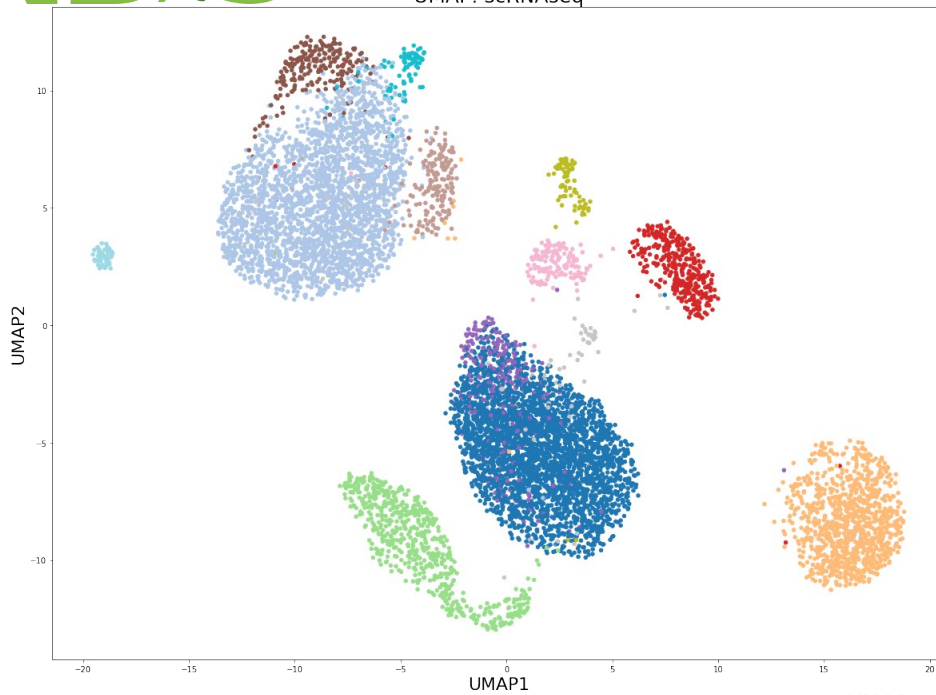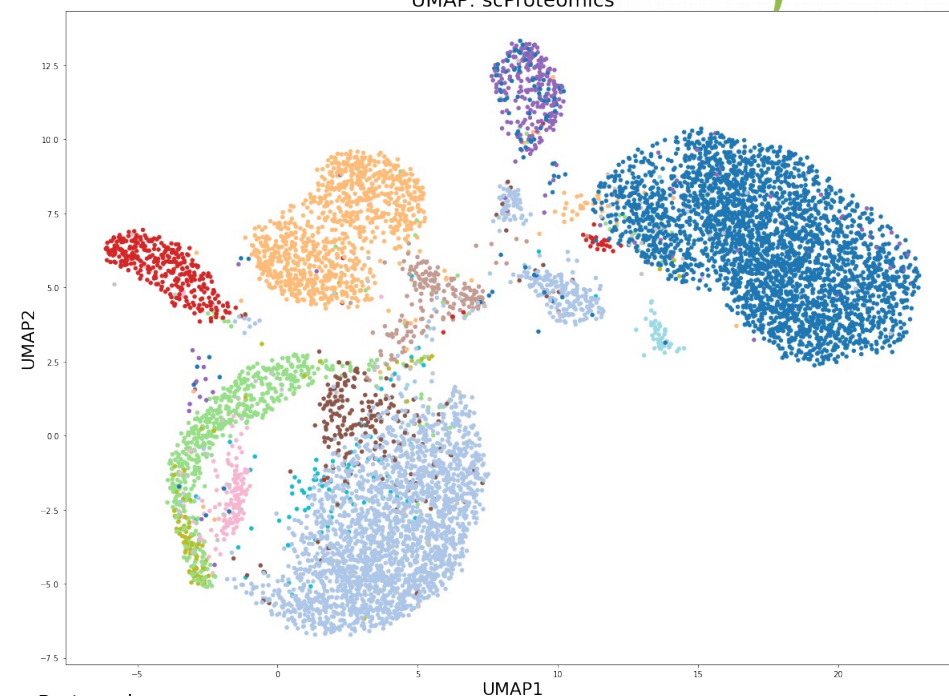
Graph Intersection

Let $S$ be a set and $F = \{S_1, \ldots, S_p\}$ a nonempty family of distinct nonempty subsets of $S$ whose union is $\bigcup_{i=1}^{p} S_i = S$. The intersection graph of $F$ is denoted $\Omega(F)$ and defined by $V(\Omega(F)) = F$, with $S_i$ and $S_j$ adjacent whenever $i \neq j$ and $S_i \cap S_j \neq \emptyset$. Then a graph $G$ is an intersection graph on $S$ if there exists a family $F$ of subsets for which $G$ and $\Omega(F)$ are isomorphic graphs (Harary 1994, p. 19). Graph intersections can be computed in the Wolfram Language using GraphIntersection[g, h].
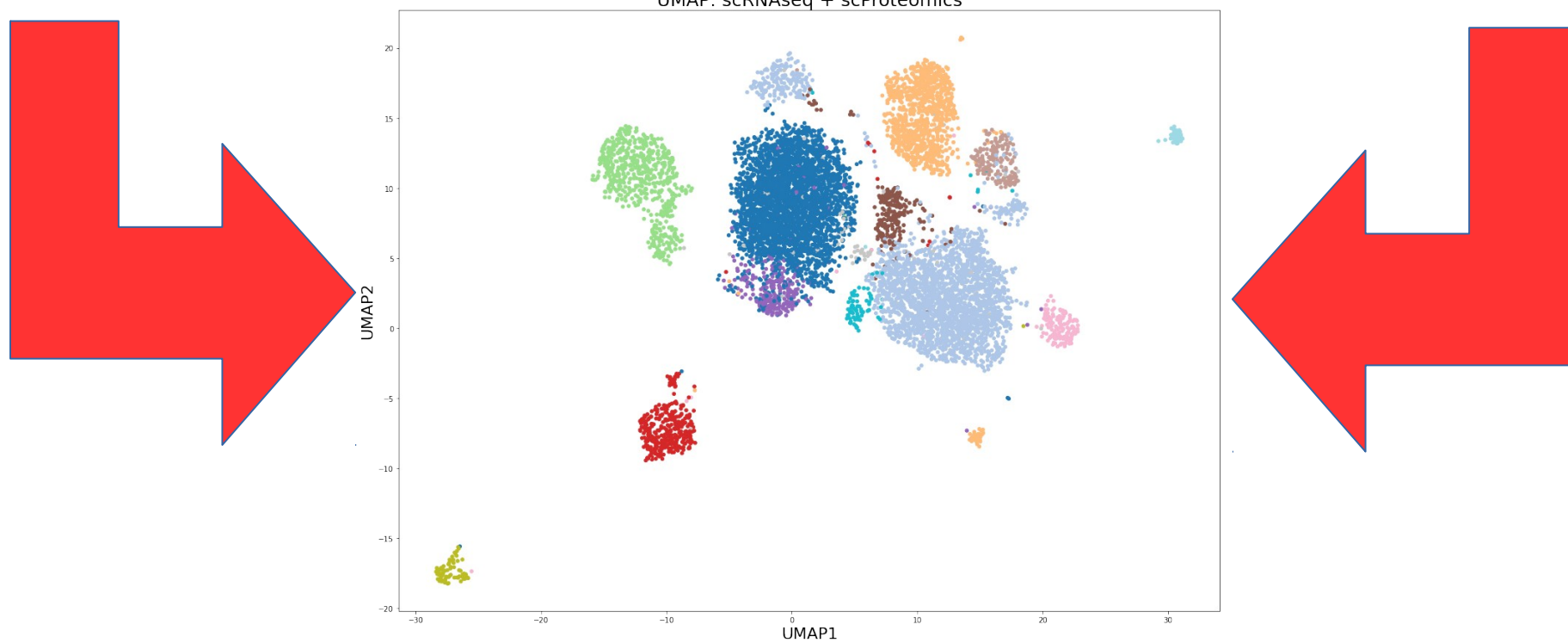
UMAP for Omics Integration
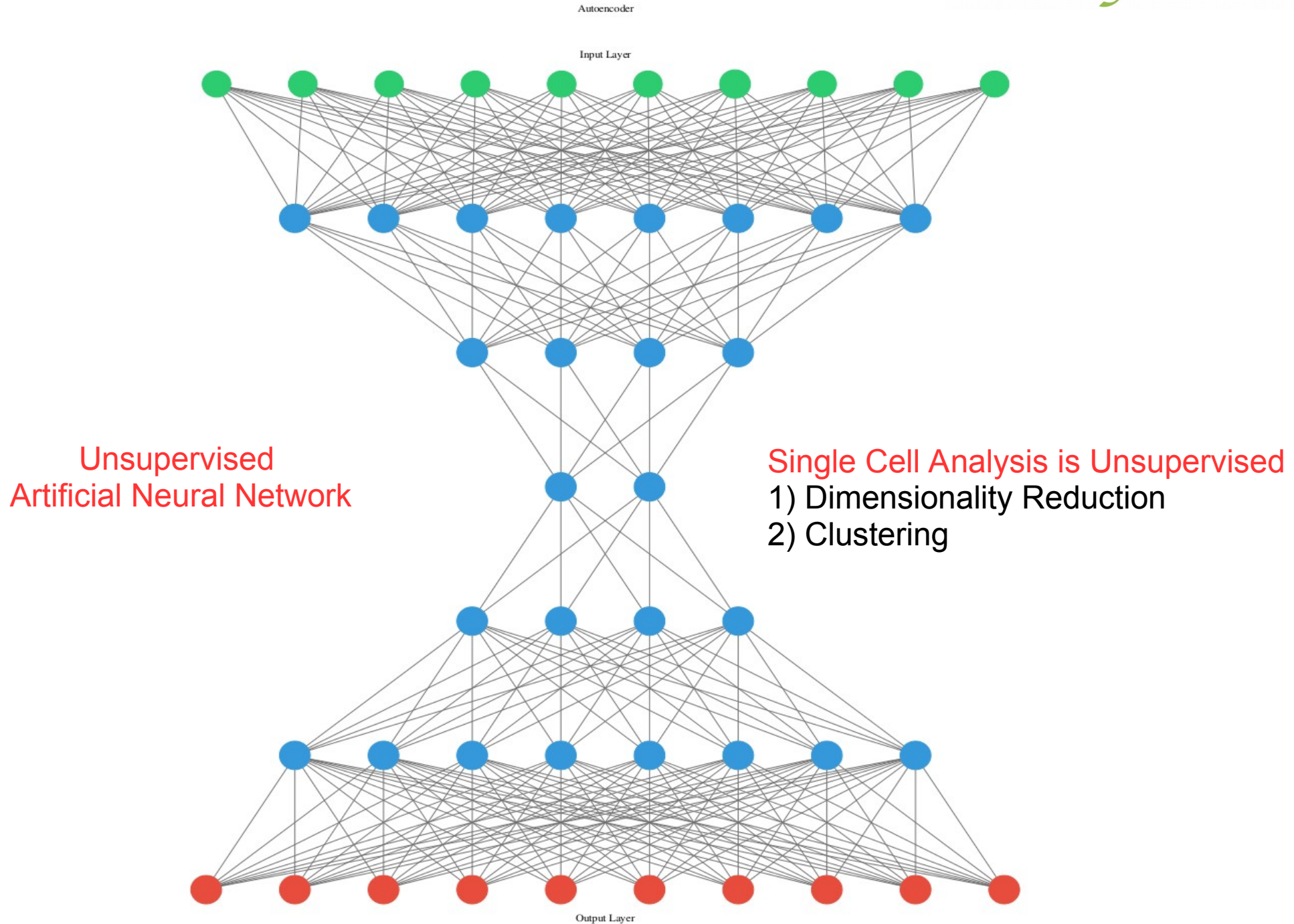
# National Bioinformatics Infrastructure Sweden (NBIS)

SciLifeLab

NBIS

Knut och Alice Wallenbergs Stiftelse

Vetenskapsrådet

LUNDS UNIVERSITET