

Machine learning methods for single-cell transcriptomic data

Guido Sanguinetti

Data Science @ SISSA, Trieste

Artificial Intelligence and Modern Physics:
a two-way Connection
Monopoli, Oct 2024

Roadmap of today

- 1 The questions and the data
- 2 Problems and approaches in single-cell 'omics
- 3 NeuroVelo: dynamics from scRNA-seq
- 4 Conclusions and perspectives

Roadmap of today

- 1 The questions and the data
- 2 Problems and approaches in single-cell 'omics
- 3 NeuroVelo: dynamics from scRNA-seq
- 4 Conclusions and perspectives

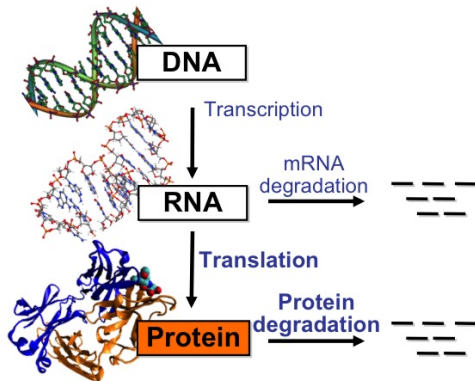
And now for something completely different ...

- In physics, we have usually a good characterisation of observables/ state variables
- We also have good ideas about the structure of the models that describe the interactions

And now for something completely different ...

- In physics, we have usually a good characterisation of observables/ state variables
- We also have good ideas about the structure of the models that describe the interactions
- Biomedicine is different: it has all the hallmarks of complex systems, but almost no theories...

The central dogma



Where does variability come to play?

Epigenetics

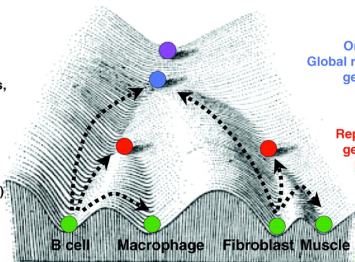
Developmental potential

Totipotent
Zygote

Pluripotent
ICM/ES cells, EG cells,
EC cells, mGS cells
iPS cells

Multipotent
Adult stem cells
(partially
reprogrammed cells?)

Unipotent
Differentiated cell
types



Epigenetic status

Global DNA demethylation

Only active X chromosomes;
Global repression of differentiation
genes by Polycomb proteins;
Promoter hypomethylation

X inactivation;
Repression of lineage-specific
genes by Polycomb proteins;
Promoter hypermethylation

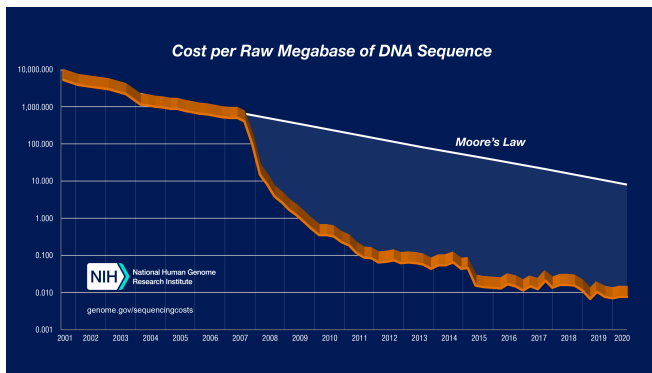
X inactivation;
Derepression of
Polycomb silenced
lineage genes;
Promoter hypermethylation

There must be something outside of (epi) the genetics that explains the variety of cell states. How to measure a cell's state?

What is sequencing?

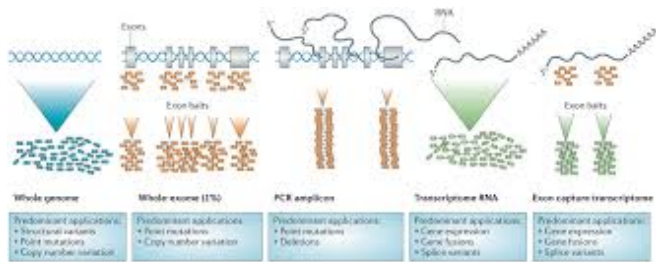
- Sequencing: (chemical or otherwise) procedure for reading the sequence of constituents of a stretch of DNA (or a protein)
- IMPORTANT 1: all DNA sequencing happens in short chunks (from a few tens to a few tens of thousands bp)
- IMPORTANT 2: mistakes are possible (depending on the technology) but quality controls available (Phred scores)
- IMPORTANT 3: some sequences are easier to sequence (biases in the data)

More than Moore



Source: National Human Genome Research Institute

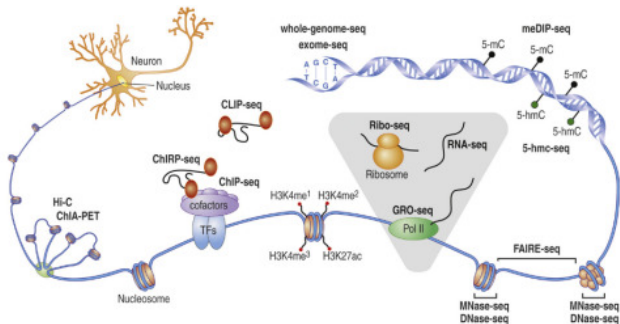
Now sequencing is easier?



Nature Reviews | Drug Discovery

Major technical advances + massive parallelisation, technology started coming online around 2008. Throughput (most recent versions) in the region 10^5 Mb/hr.

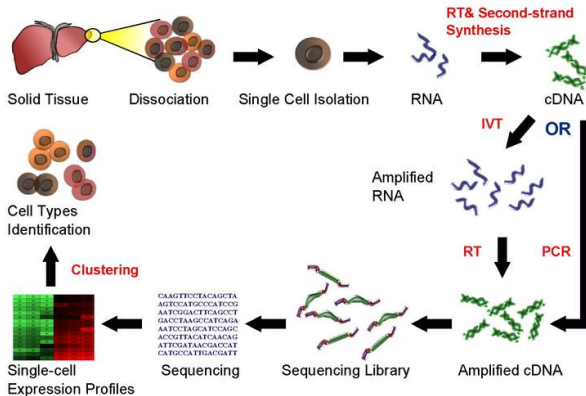
Sequencing everywhere



- NGS can be "prefaced" by any biochemical treatment
- IMPORTANT: when doing that biases are often introduced/
becomes unclear how to compare samples

scRNA-seq

Single Cell RNA Sequencing Workflow



Can do 100K cells in single experiment. High dropout rate, huge variability in coverage. Dominant technology now.

Roadmap of today

- 1 The questions and the data
- 2 Problems and approaches in single-cell 'omics**
- 3 NeuroVelo: dynamics from scRNA-seq
- 4 Conclusions and perspectives

What single-cell 'omics look like

- For each cell, we normally obtain $\sim 10\text{K}$ RNA fragments mapped to the transcriptome \rightarrow *most genes* are missed in every single cell
- We apply some pre-filtering criterion, e.g. discard genes not measured in at least 50% of cells, cells with fewer than 100 non-zero genes
- We end up with a gene expression matrix typically $\sim 6\text{K}$ rows (genes) and a few thousands columns (cells)

What single-cell 'omics look like

- For each cell, we normally obtain $\sim 10\text{K}$ RNA fragments mapped to the transcriptome \rightarrow *most genes* are missed in every single cell
- We apply some pre-filtering criterion, e.g. discard genes not measured in at least 50% of cells, cells with fewer than 100 non-zero genes
- We end up with a gene expression matrix typically $\sim 6\text{K}$ rows (genes) and a few thousands columns (cells)
- A large fraction of the entries are zero, either genuine or dropout

The general problem of Data Science

- We have noisy, high dimensional data
- Several factors contribute to the variance we observe in the data: experimental noise, intrinsic stochasticity, physiological processes, disease processes

The general problem of Data Science

- We have noisy, high dimensional data
- Several factors contribute to the variance we observe in the data: experimental noise, intrinsic stochasticity, physiological processes, disease processes
- **ALL** of Data Science/ AI consists in partitioning this variance and using this (implicit or explicit) decomposition for predictions

The general problem of Data Science

- We have noisy, high dimensional data
- Several factors contribute to the variance we observe in the data: experimental noise, intrinsic stochasticity, physiological processes, disease processes
- **ALL** of Data Science/ AI consists in partitioning this variance and using this (implicit or explicit) decomposition for predictions
- The difference lies in the assumptions about what is an important direction of variation

The graph abstraction

- One general approach to deal with noisy data is to threshold features
- When considering a system, this can be conceptualised as setting to zero some interactions \rightarrow building a graph

The graph abstraction

- One general approach to deal with noisy data is to threshold features
- When considering a system, this can be conceptualised as setting to zero some interactions \rightarrow building a graph
- Many popular methods in the analysis of scRNA-seq are based on a graph abstraction

Graphs of cells

- Here the entities are cells, and the feature we measure high-dimensional expression readouts

Graphs of cells

- Here the entities are cells, and the feature we measure high-dimensional expression readouts
- Using some procedure (e.g. PCA followed by Pearson correlation) we define a (dis-) similarity measure between cells

Graphs of cells

- Here the entities are cells, and the feature we measure high-dimensional expression readouts
- Using some procedure (e.g. PCA followed by Pearson correlation) we define a (dis-) similarity measure between cells
- Then each cell is connected only with the K most similar cells (K -nearest neighbour graph) or with cells with a similarity that exceeds a certain threshold
- The edges are typically weighed by the similarity

Graphs of cells

- Here the entities are cells, and the feature we measure high-dimensional expression readouts
- Using some procedure (e.g. PCA followed by Pearson correlation) we define a (dis-) similarity measure between cells
- Then each cell is connected only with the K most similar cells (K -nearest neighbour graph) or with cells with a similarity that exceeds a certain threshold
- The edges are typically weighed by the similarity
- Subsequent results may (and do) depend on the hyperparameters (thresholds/ K / pre-filtering)

Problem 1: Visualisation/ dimensionality reduction

- Assumption: Only few degrees of freedom exist in the data

Problem 1: Visualisation/ dimensionality reduction

- Assumption: Only few degrees of freedom exist in the data
- Dominant tool: UMAP (Uniform Manifold Approximation (McInnes et al 2018))
- Creates a nearest neighbour graph in gene space, then tries to find points in low D such that the graph distances are preserved

Problem 1: Visualisation/ dimensionality reduction

- Assumption: Only few degrees of freedom exist in the data
- Dominant tool: UMAP (Uniform Manifold Approximation (McInnes et al 2018))
- Creates a nearest neighbour graph in gene space, then tries to find points in low D such that the graph distances are preserved
- Because it tries to inherit the geometric properties of the high dimensional graph, it may show structures that are not obviously there
- Not easy to understand what UMAP directions mean

Sub-problem 2: Pseudo-time

- Assumption: the major direction of variation is along a developmental direction
- E.g., cells are collectively following a dynamical process (development, differentiation, drug response) but individually they are at slightly different stages of the process

Sub-problem 2: Pseudo-time

- Assumption: the major direction of variation is along a developmental direction
- E.g., cells are collectively following a dynamical process (development, differentiation, drug response) but individually they are at slightly different stages of the process
- Given a root cell, the pseudotime of a cell becomes the expected time to reach it from the root by random walk on the graph (Diffusion Maps, Hagheverdi 2015)
- Output: (partial) ordering of cells, identification of branching events

Problem 3: clustering

- Assumption: the major variation is caused by the existence of distinct groups of cells which are transcriptomically homogeneous
- Solution: identify clusters or communities (densely connected regions of the graph)
- Popular ones based on optimising the *modularity* of the graph w.r.t. a partitioning of the graph
- Most commonly used algorithms are the Louvain algorithm (Blondel et al 2008) and the Leiden algorithm (Traag et al 2019)

Software packages



Q Search

Installation

Tutorials

Usage Principles

How to

API

Stars 1.8k pypi v1.10.1 downloads 3M downloads 1.39k docs passing Azure Pipelines never built discourse 4.3k posts zulip join chat powered by NumFOCUS

Scanpy – Single-Cell Analysis in Python

Scanpy is a scalable toolkit for analyzing single-cell gene expression data built jointly with [anndata](#). It includes preprocessing, visualization, clustering, trajectory inference and differential expression testing. The Python-based implementation efficiently deals with datasets of more than one million cells.



Seurat v5

We are excited to release Seurat v5! To install, please follow the instructions in our [install page](#). This update brings the following new features and functionality:

Links

[View on CRAN](#)

[Browse source code](#)

[Report a bug](#)

License

[Full license](#)

[MIT + file LICENSE](#)

Community

[Code of conduct](#)

Both R and Python well used

Roadmap of today

- 1 The questions and the data
- 2 Problems and approaches in single-cell 'omics
- 3 NeuroVelo: dynamics from scRNA-seq**
- 4 Conclusions and perspectives

Splicing: another layer of complexity

Video of splicing

Uncovering dynamics: RNA velocity

- scRNA-seq is destructive → static snapshots from a dynamic process

Uncovering dynamics: RNA velocity

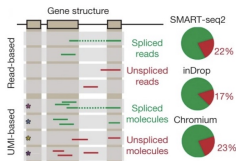
- scRNA-seq is destructive → static snapshots from a dynamic process
- **IDEA** (La Manno et al, 2018): use spliced/ unspliced reads to derive *rate of change* of RNA levels

$$\frac{dx_u}{dt} = \alpha - \beta x_u \quad \frac{dx_s}{dt} = \beta x_u - \gamma x_s$$

Uncovering dynamics: RNA velocity

- scRNA-seq is destructive → static snapshots from a dynamic process
- **IDEA** (La Manno et al, 2018): use spliced/ unspliced reads to derive *rate of change* of RNA levels

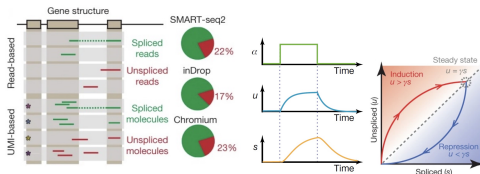
$$\frac{dx_u}{dt} = \alpha - \beta x_u \quad \frac{dx_s}{dt} = \beta x_u - \gamma x_s$$



Uncovering dynamics: RNA velocity

- scRNA-seq is destructive → static snapshots from a dynamic process
- **IDEA** (La Manno et al, 2018): use spliced/ unspliced reads to derive *rate of change* of RNA levels

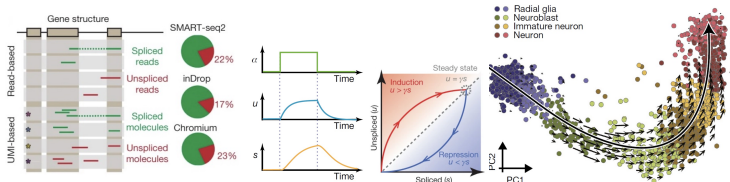
$$\frac{dx_u}{dt} = \alpha - \beta x_u \quad \frac{dx_s}{dt} = \beta x_u - \gamma x_s$$



Uncovering dynamics: RNA velocity

- scRNA-seq is destructive → static snapshots from a dynamic process
- **IDEA** (La Manno et al, 2018): use spliced/ unspliced reads to derive *rate of change* of RNA levels

$$\frac{dx_u}{dt} = \alpha - \beta x_u \quad \frac{dx_s}{dt} = \beta x_u - \gamma x_s$$



Problems and solutions

- Splicing signal is **very noisy** in single cells
- No reason why timescale of splicing should be the relevant one

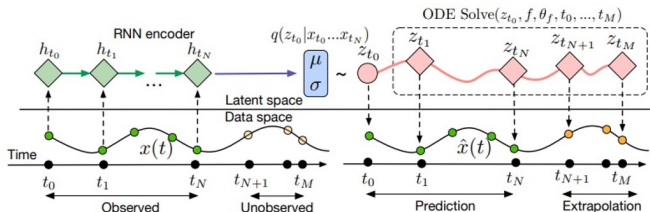
Problems and solutions

- Splicing signal is **very noisy** in single cells
- No reason why timescale of splicing should be the relevant one
- **IDEA:** Underlying (low dimensional) nonlinear dynamical system should govern long-term evolution of cells' transcriptomes
- Spliced/ unspliced ratio gives a noisy measurement of *instantaneous* rate of change

Problems and solutions

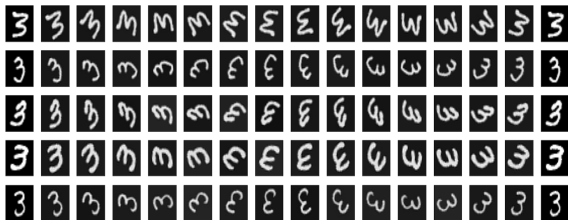
- Splicing signal is **very noisy** in single cells
- No reason why timescale of splicing should be the relevant one
- **IDEA:** Underlying (low dimensional) nonlinear dynamical system should govern long-term evolution of cells' transcriptomes
- Spliced/ unspliced ratio gives a noisy measurement of *instantaneous* rate of change
- Couple the two components in the spirit of *physics informed machine learning*

Neural ODEs (Chen et al 2018)

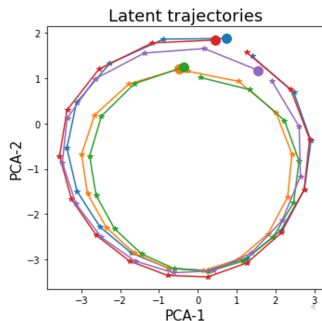
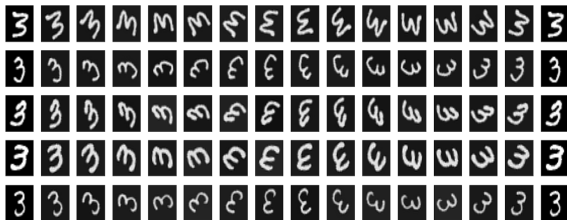


Autoencoding structure in time. ODE in latent space with drift parametrised by a NN. Efficient evaluation of gradients by Pontryagin adjoint.

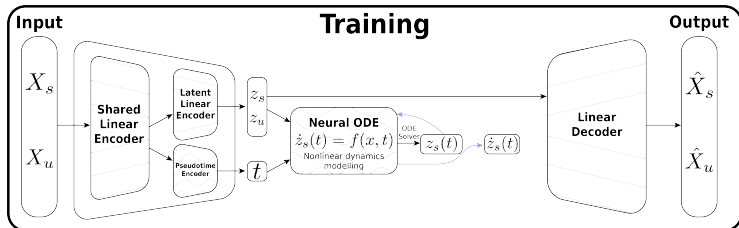
A (somewhat contrived) example



A (somewhat contrived) example



NeuroVelo (Idris Kouadri Boudjelthia)



$$\mathcal{L} = \text{MSE}(X, \hat{X}) + \text{MSE}(\dot{z}_s, \beta z_u - \gamma z_s)$$

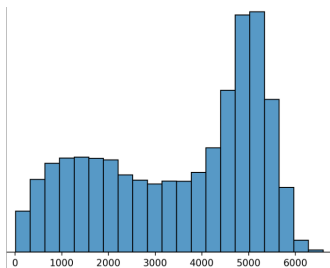
Because the encoding/ decoding is linear, the RNA velocity equations apply also in latent space. Notice we need no assumptions on the transcription rate function.

Interpreting Neurovelo

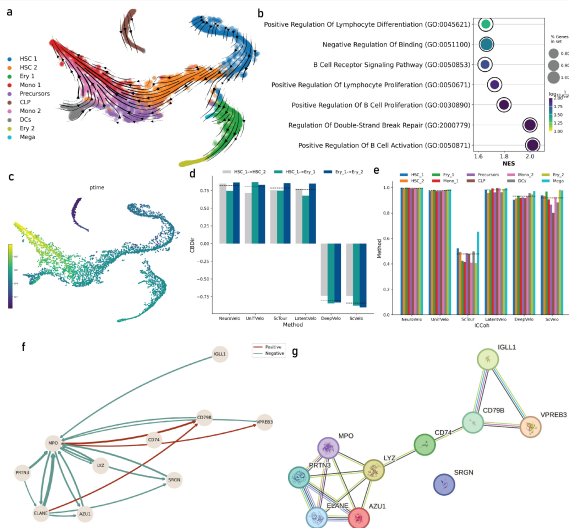
- NeuroVelo learns a low-dimensional nonlinear dynamical system
- Principal dynamics are given (locally) by the *eigenvectors* of the Jacobian matrix
- These eigenvectors can be decoded linearly to give a ranked list of genes
- The decoded Jacobian matrix gives itself a description of the network of interactions between genes
- Robustness is ensured by computing a stability index w.r.t. multiple initializations

Interpreting Neurovelo cont'd

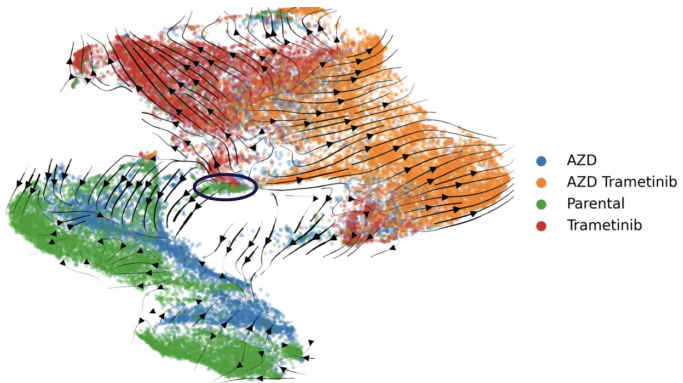
- Noise genes should have ranks uniformly distributed
← Gaussian average
- Relevant genes should have consistently high ranks
- Expect bimodal distribution



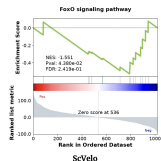
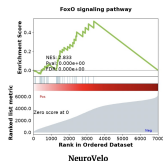
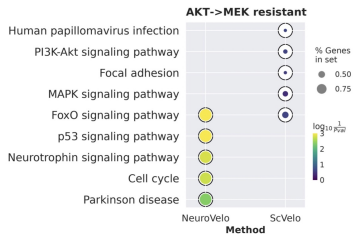
NeuroVelo on HBM



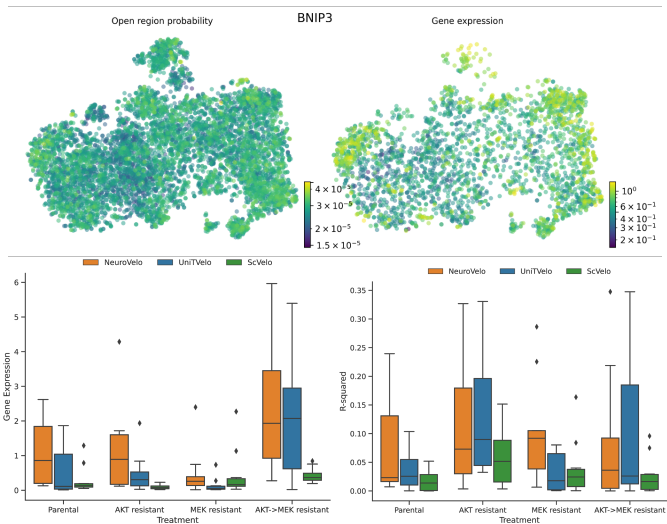
NeuroVelo on CRC



Validating NeuroVelo: enrichment



Validating NeuroVelo: multiome



Roadmap of today

- 1 The questions and the data
- 2 Problems and approaches in single-cell 'omics
- 3 NeuroVelo: dynamics from scRNA-seq
- 4 Conclusions and perspectives**

Conclusions

- Single-cell 'omics provide a potential goldmine, but you need the right pick-axe
- Must go beyond simply plotting cells in latent space
- Combining interpretability and nonlinearity is still a major challenge
- Interpretability is key to progress to the clinic!

Thanks!

Collaborators and lab members/ alumni

SISSA

Riccardo Margiotta

Rongrong Xie

Viplove Arora

Matteo Santoro

Nour el Kazwini

Alex Zhang

Idris Kouadri

Boudjelthia

Federico Caretti

Katsiaryna

Davydzenka

**University of
Edinburgh**

Kashyap Chhatbar

Kaan Ocal

Christos Maniatis

Andreas Kapourani

Yuanhua Huang

Catalina Vallejos

**Human
Technopole**

Andrea Sottoriva

Salvatore Milite

Clara Canavese

Funding: ERC, AIRC, SISSA/ MUR

Some references

- Kouadri-Boudjelthia, Idris, et al <https://www.biorxiv.org/content/10.1101/2023.11.17.567500v1>
- La Manno, Gioele et al, *Nature* 460 (494-498) (2018), <https://www.nature.com/articles/s41586-018-0414-6>
- McInnes, L. et al, *J. Open Source Softw.* 3(29) (2018) <https://joss.theoj.org/papers/10.21105/joss.00861>
- Haghverdi, L. et al, *Nature Meth.* 13 (845-848) (2016), <https://www.nature.com/articles/nmeth.3971>
- Blondel, V et al *J. Stat. Mech.: Th. Exp.* 2008 (10): 10008 <https://iopscience.iop.org/article/10.1088/1742-5468/2008/10/P10008>
- Traag, V et al *Scientific Reports.* 9 (1): 5233. <https://arxiv.org/abs/1810.08473>