



# Topological Data Analysis: extracting insights from the “shape” of data

55e Journées de statistique de la SFdS  
27 May 2024

**EPFL**

What is topology?

A word cloud centered around the word "topology". The word "topology" is the largest and most prominent, written in a light blue color. Other large words include "shape" (white, vertical), "connectivity" (blue, horizontal), "algebra" (blue, horizontal), and "geometry" (yellow, vertical). Smaller words include "deformation" (orange), "path" (blue), "classification" (orange), "invariants" (orange), "cavity" (orange), "open" (blue), "homology" (white), "equivalence" (orange, vertical), "closed" (blue), "continuity" (blue), "complex" (blue), "mug" (orange), "connected" (orange), "simplex" (orange), "donut" (white), "proximity" (blue), and "continuous" (blue).

deformation  
connectivity  
shape  
continuous  
path  
cavity  
invariants  
donut  
simplex  
classification  
open  
topology  
homology  
algebra  
equivalence  
closed  
continuity  
complex  
mug  
connected  
geometry

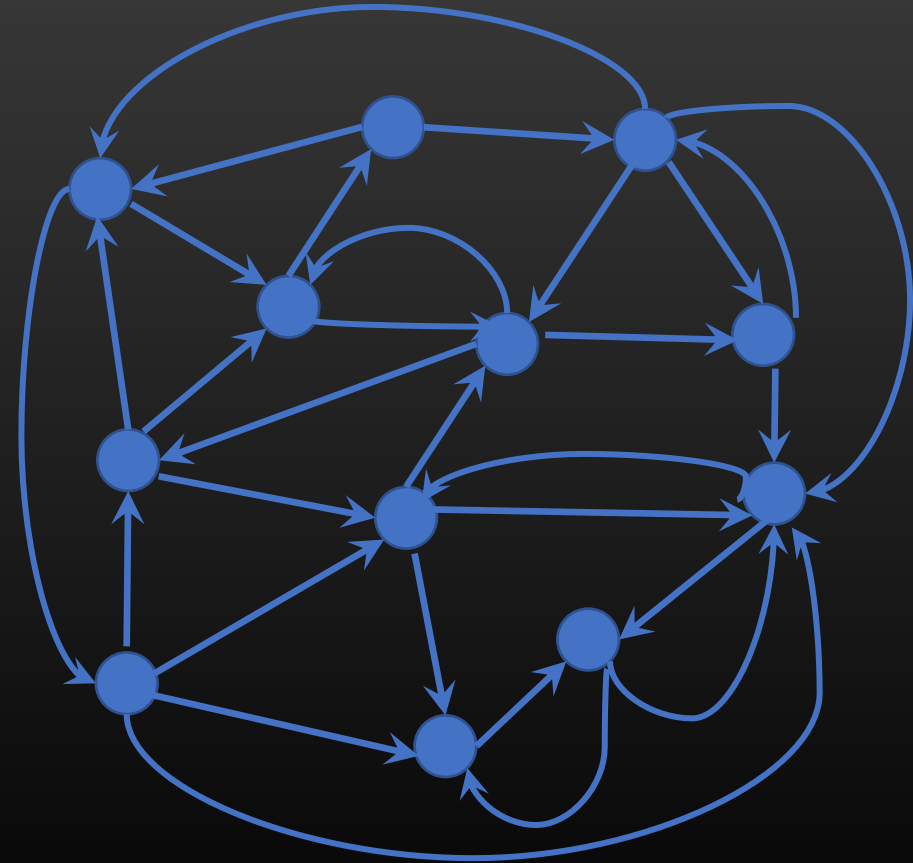
# Topology is...

- the mathematics of **shape**;



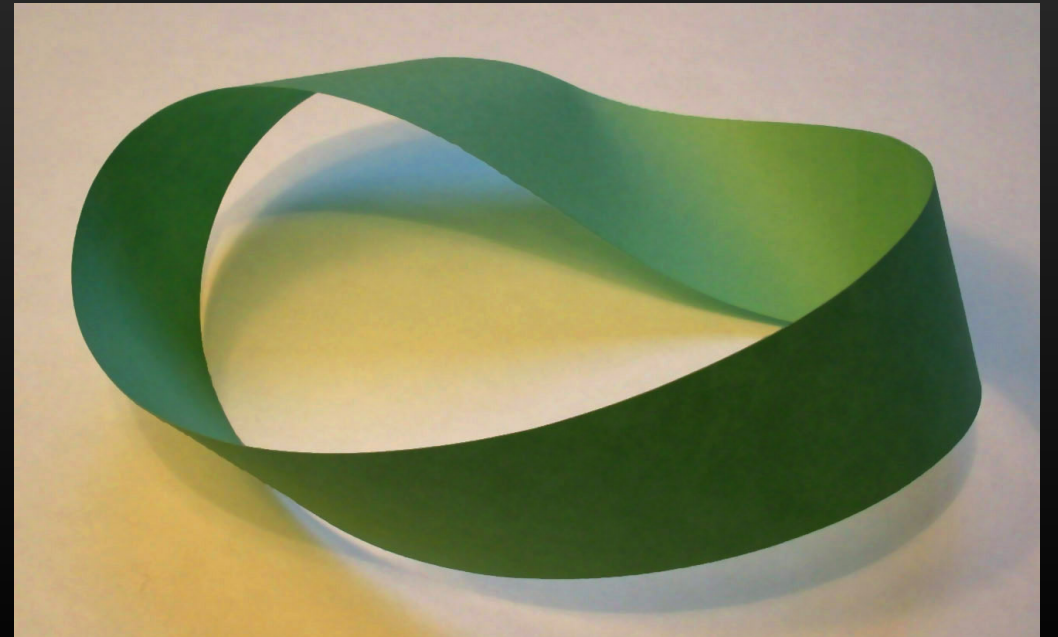
# Topology is...

- the mathematics of **shape**;
- the mathematics of **connectivity**;



# Topology is...

- the mathematics of **shape**;
- the mathematics of **connectivity**;
- the mathematics of **emergence of global structure from local constraints**.



# Application to Data Science

The **shape** of a data set,  
described by a **topological signature**  
encoding its **multi-scale structure**,  
can reveal important relations among the data points,  
with the help of machine learning.

**Topological Data Analysis (TDA)**

# Topological analytical tools

Method	Appropriate data types
Mapper	Clinical data, metabolomics, genomics, etc.
Two-tier Mapper	Gene expression data, single-cell transcriptomics
Persistent homology	Connectivity data, high-dimensional point cloud data
Graph signal processing	Connectivity data + “signal”

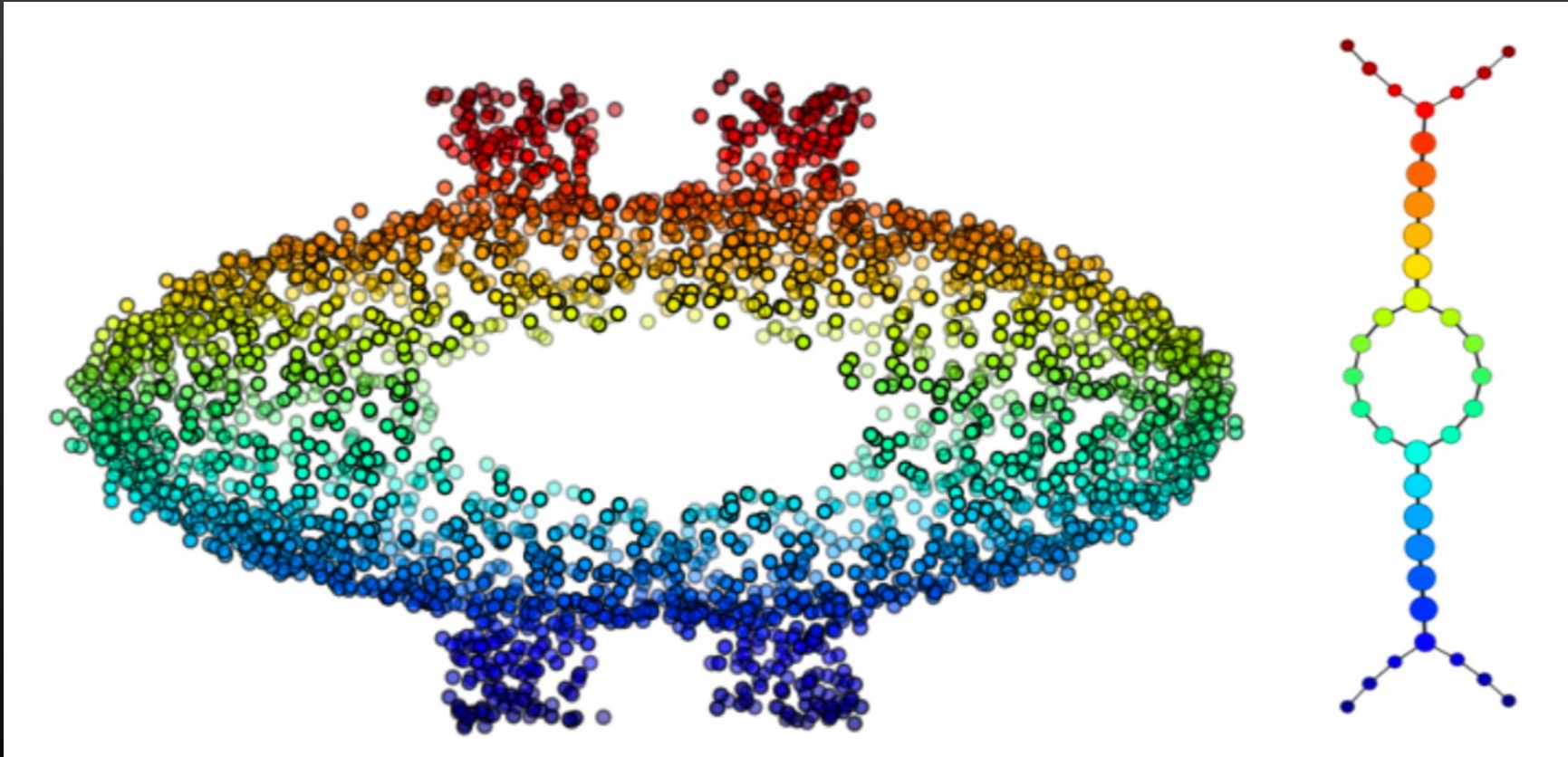


Mapper

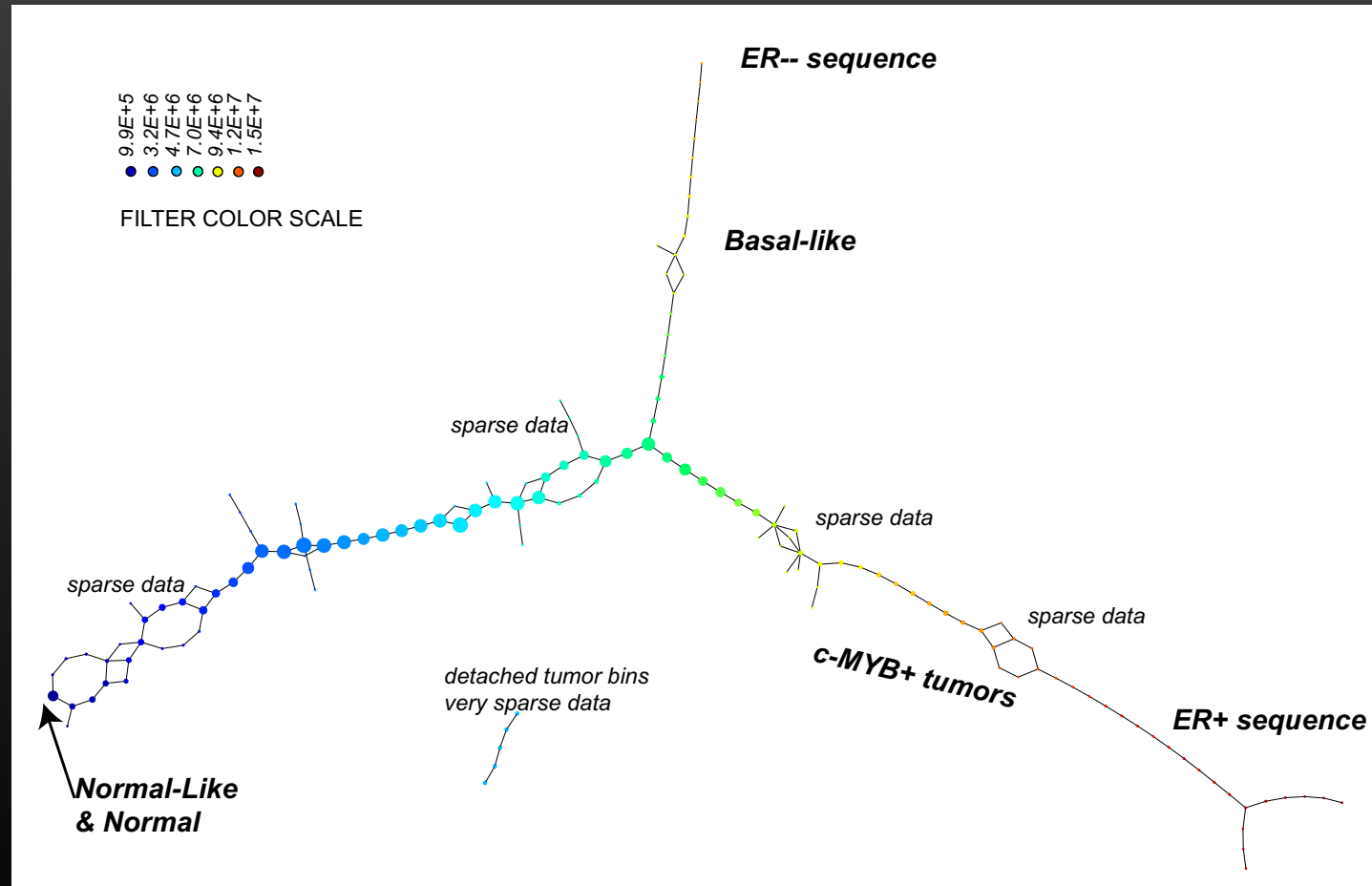
# Overview

- (Mostly) **unsupervised** multivariate pattern analysis of high-dimensional data, retains **more information than PCA**
- Produces a **compressed visual representation** of the data, providing a strong indication of where to look for **meaningful clustering** and encoding **relations between clusters**
- Numerous remarkably successful applications
- Input:
  - Data set  $X$  equipped with notion of “distance” between points
  - Real-valued “measurements” on  $X$
  - Decomposition of the real line into overlapping subsets
  - Choice of clustering algorithm

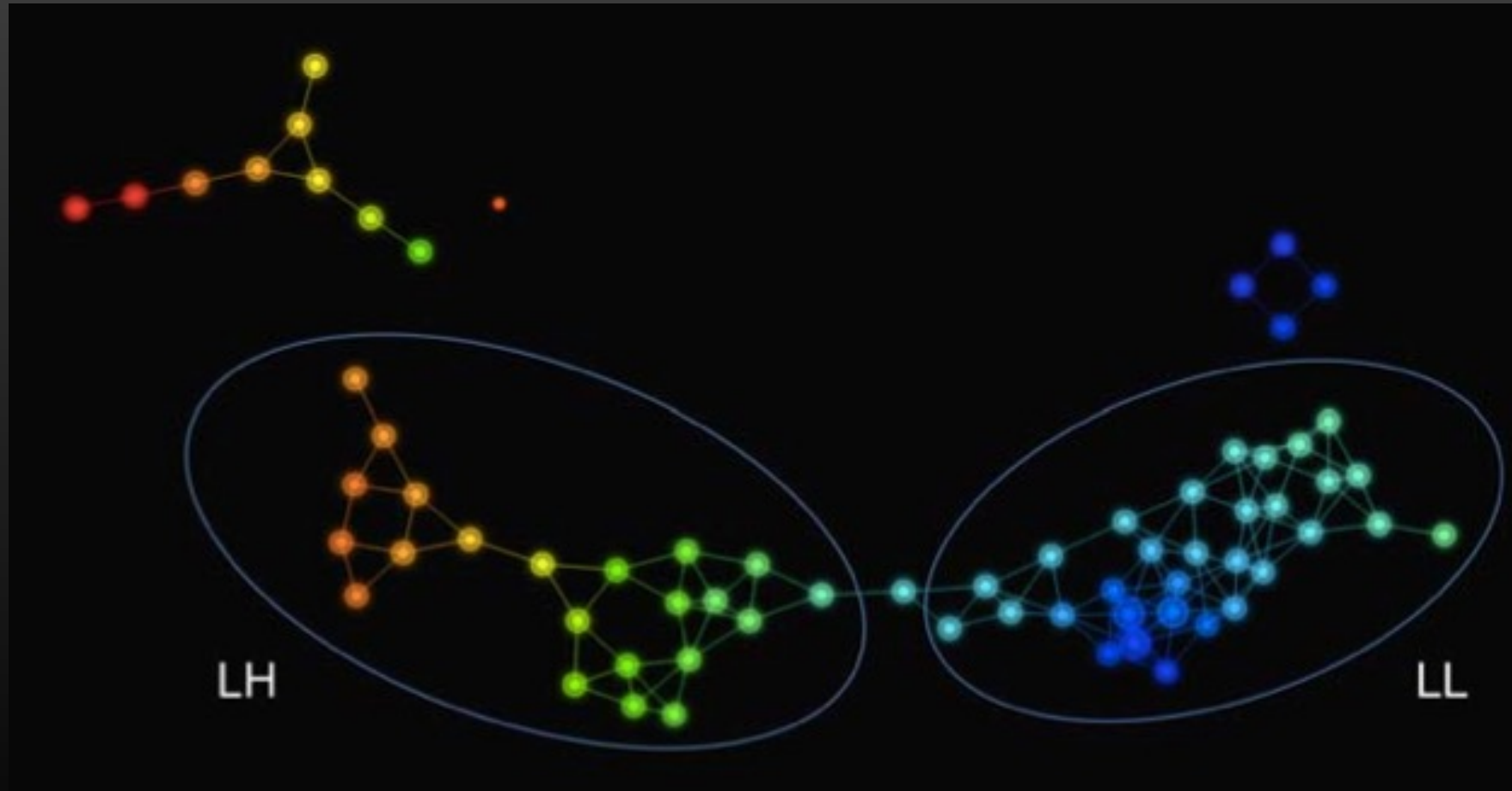
# Mapper output: synthetic data



# Mapper output: gene expression data



# Mapper output: fMRI data



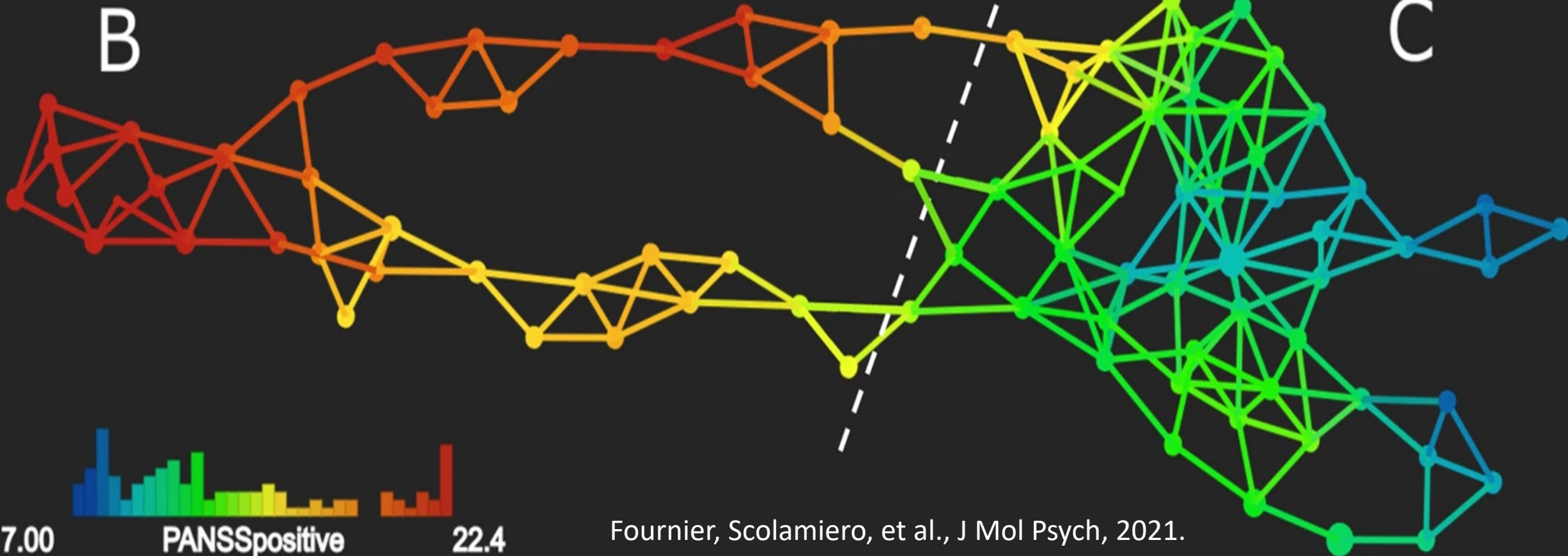
A



Mapper output: clinical profile data



B



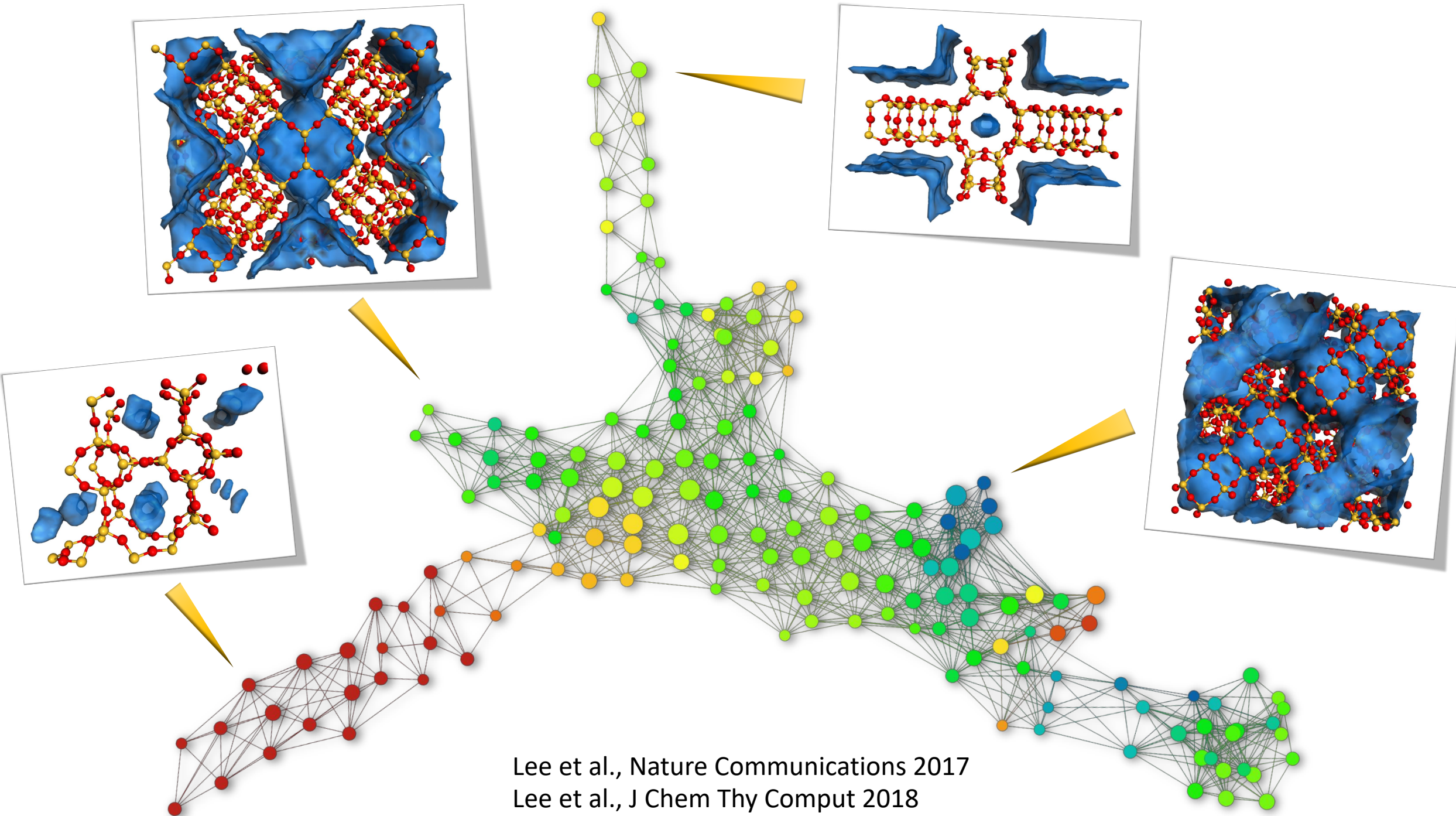
C

7.00

PANSSpositive

22.4

Fournier, Scolamiero, et al., J Mol Psych, 2021.

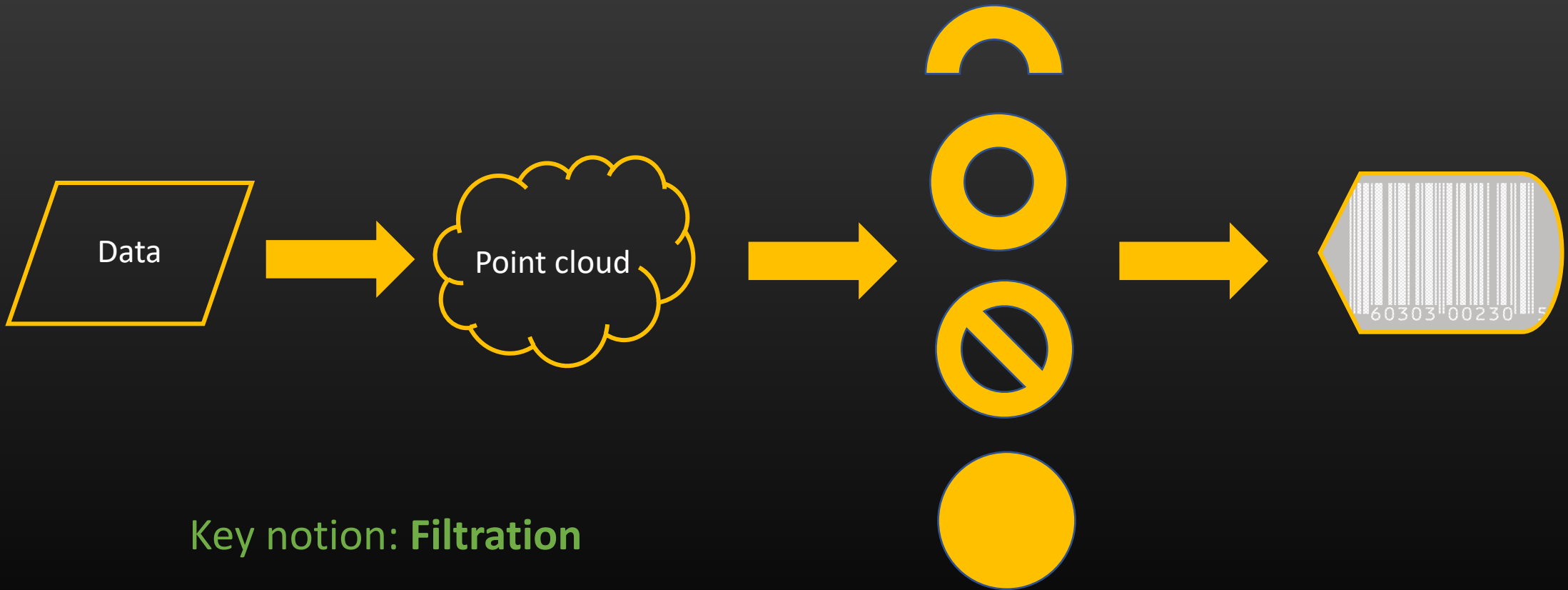


Lee et al., Nature Communications 2017  
Lee et al., J Chem Thy Comput 2018

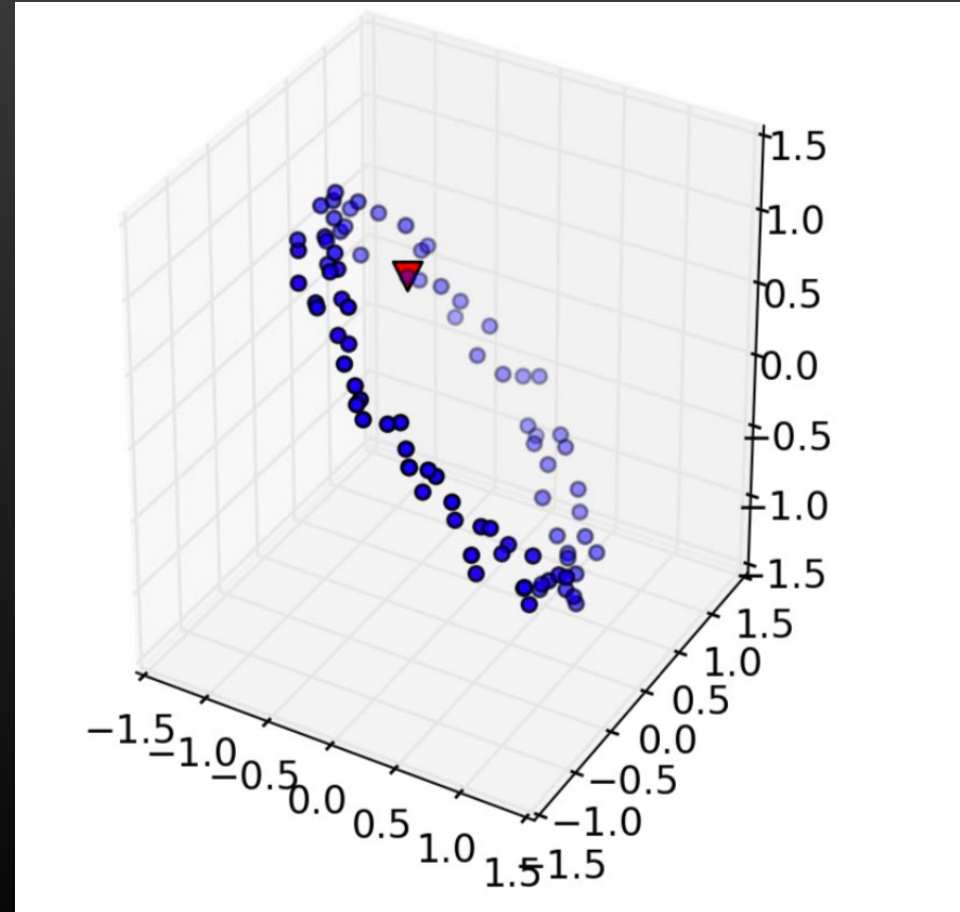
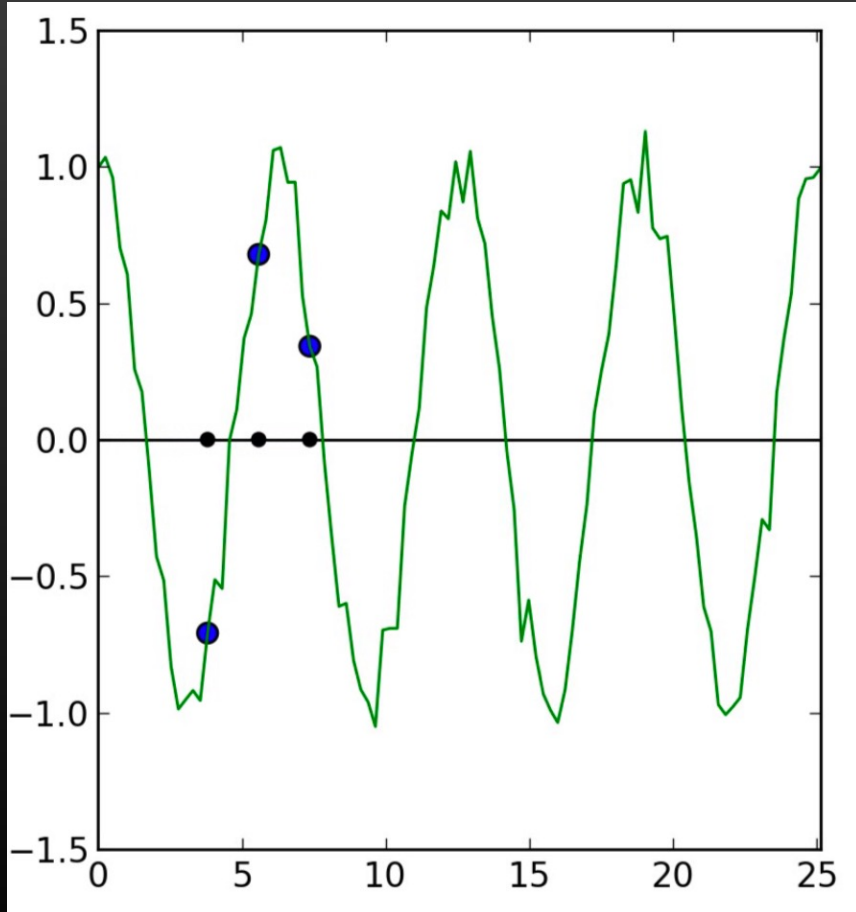
Persistent homology



# The basic persistence workflow

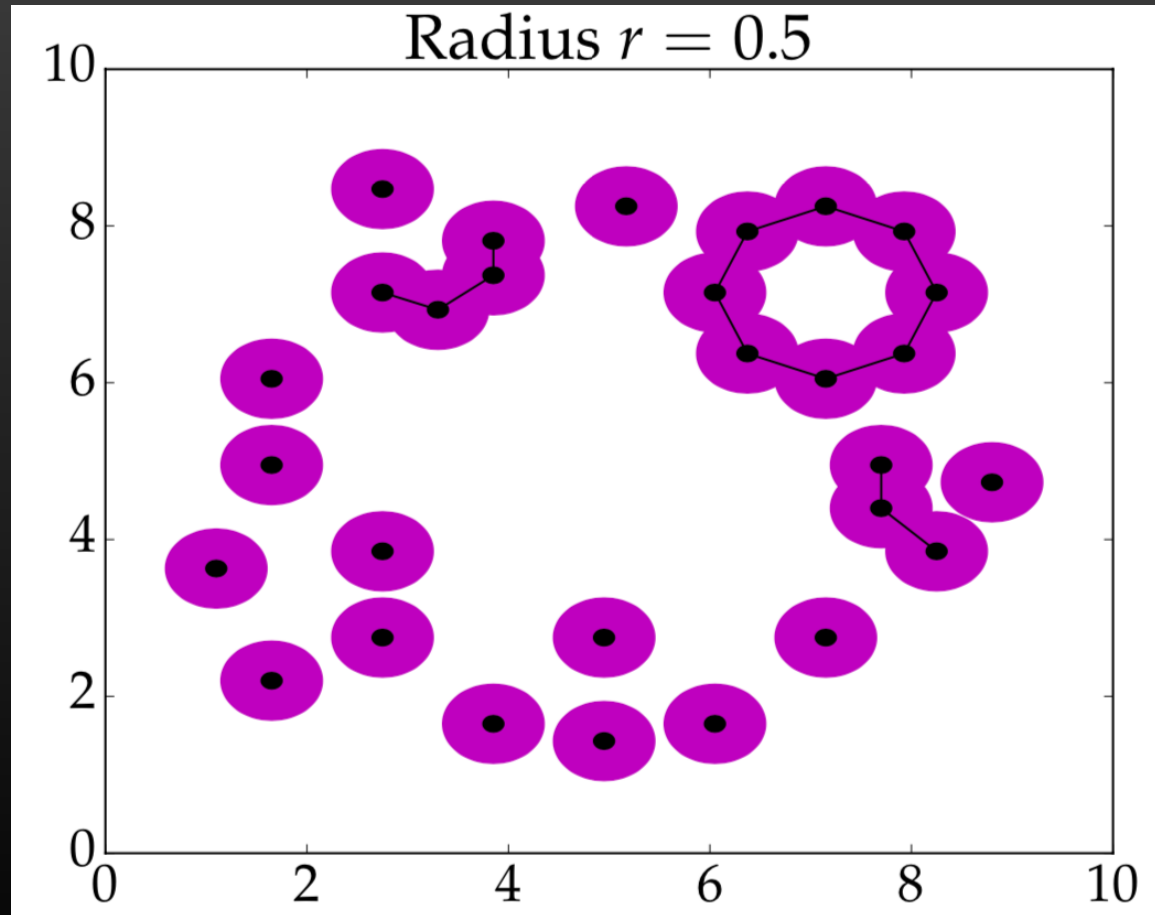


# Step 1: Data to Point Cloud

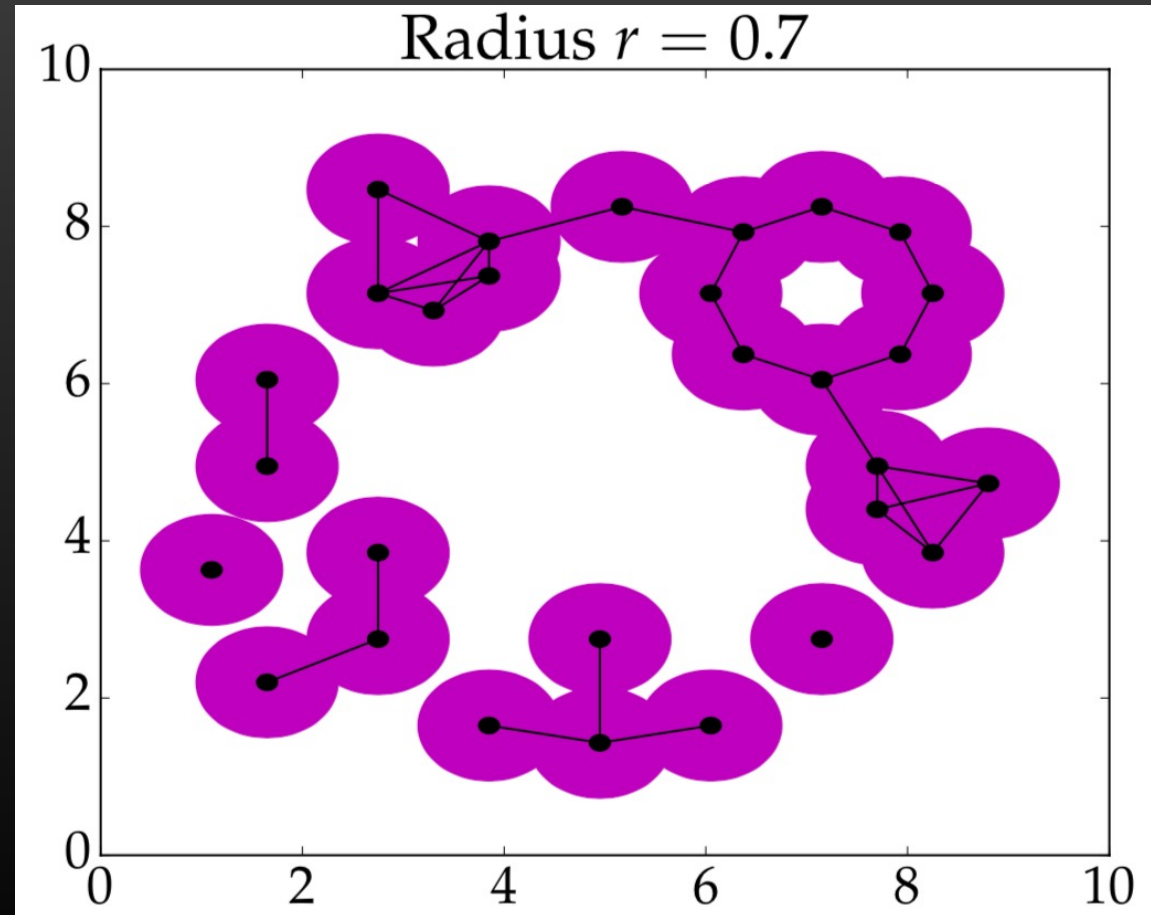




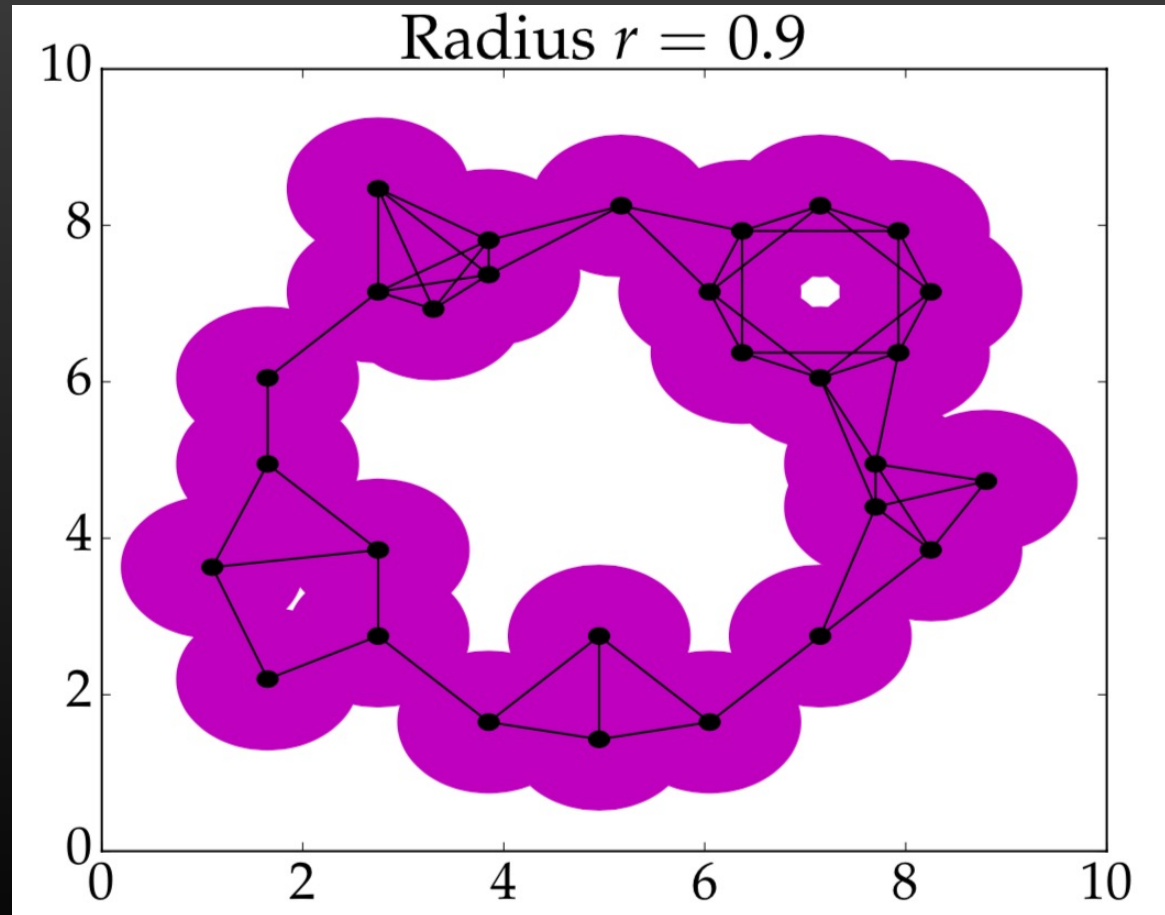
# Step 2: Point cloud to nested complexes



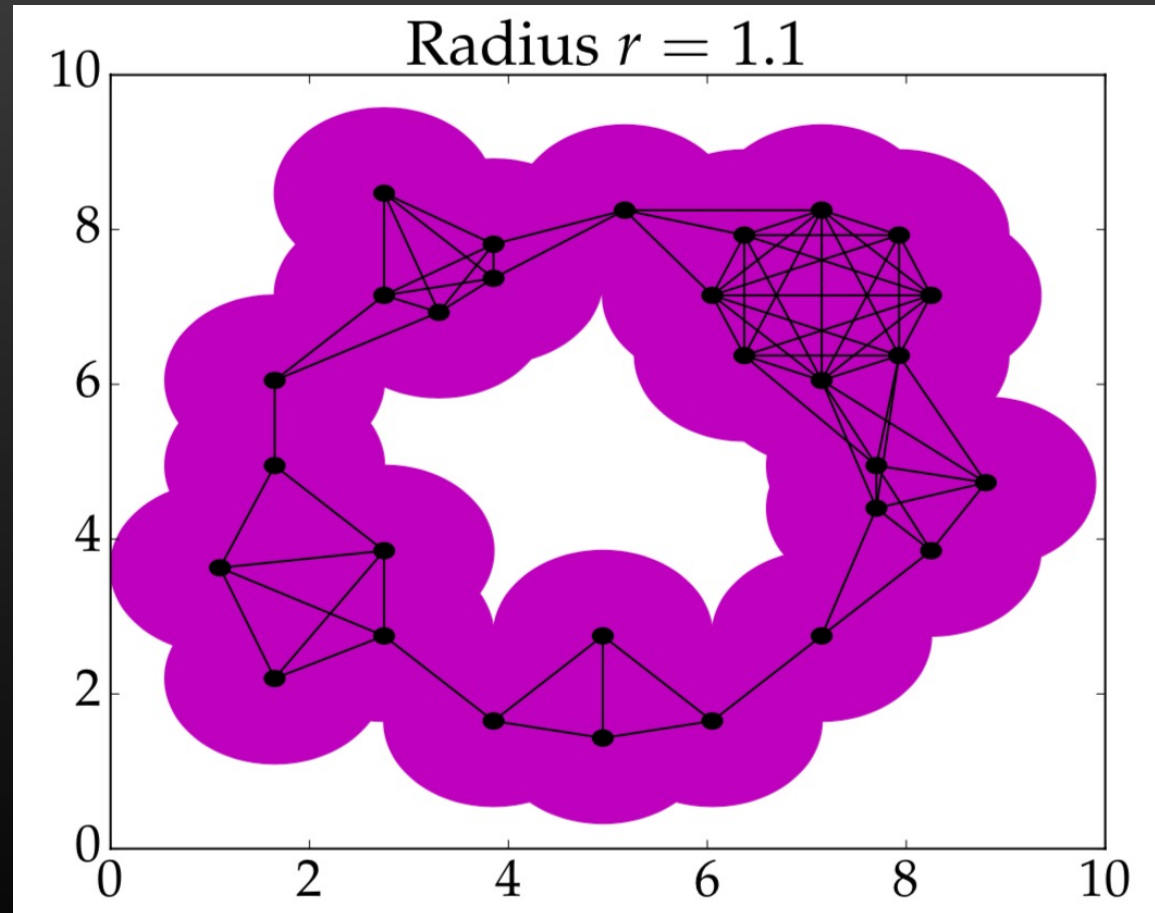
# Step 2: Point cloud to nested complexes



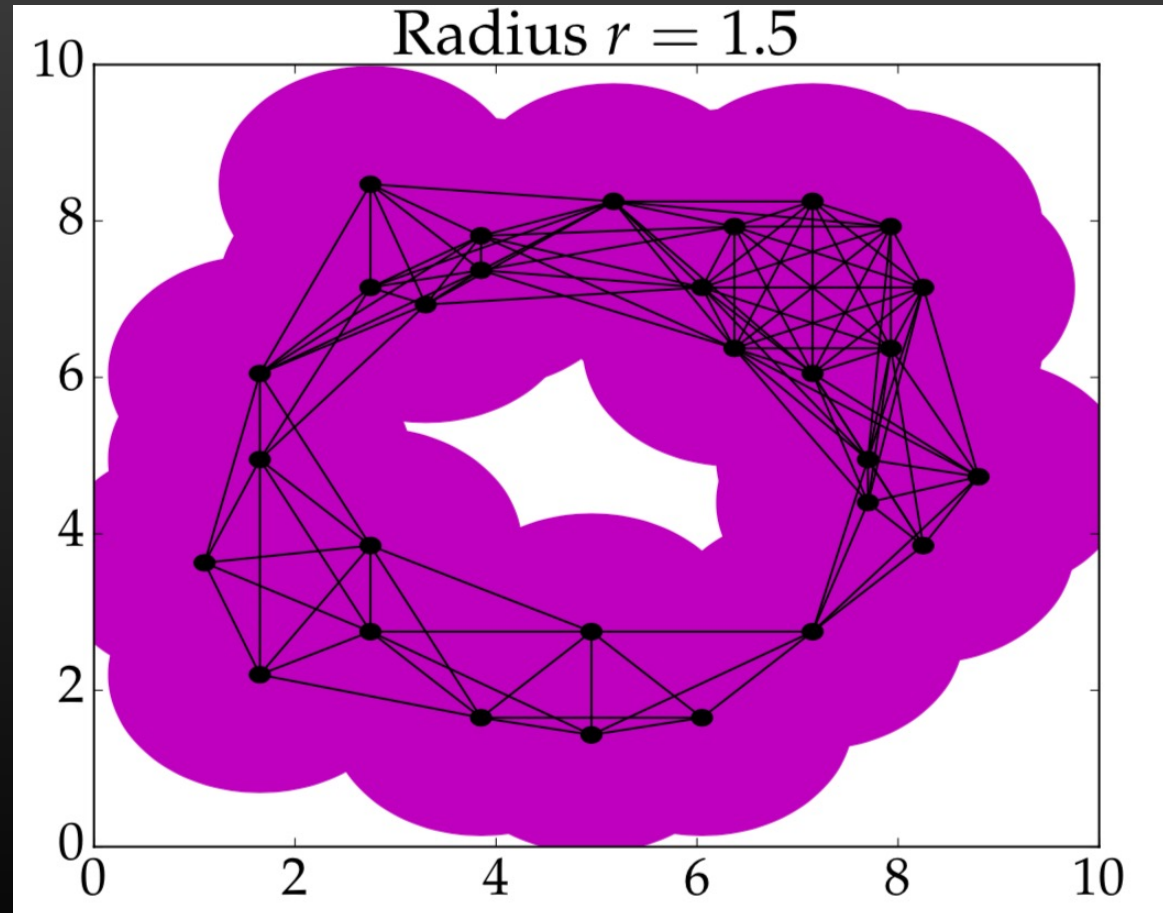
# Step 2: Point cloud to nested complexes



# Step 2: Point cloud to nested complexes

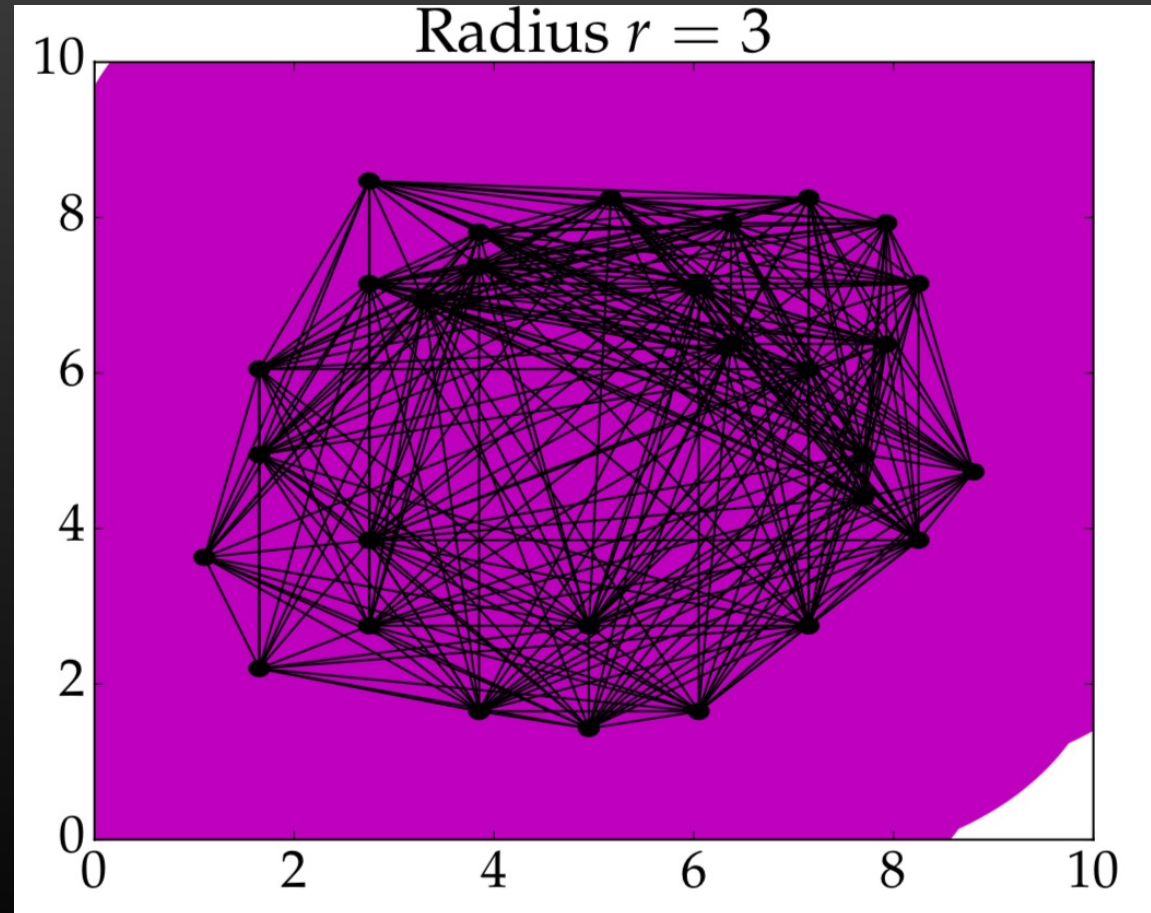


# Step 2: Point cloud to nested complexes

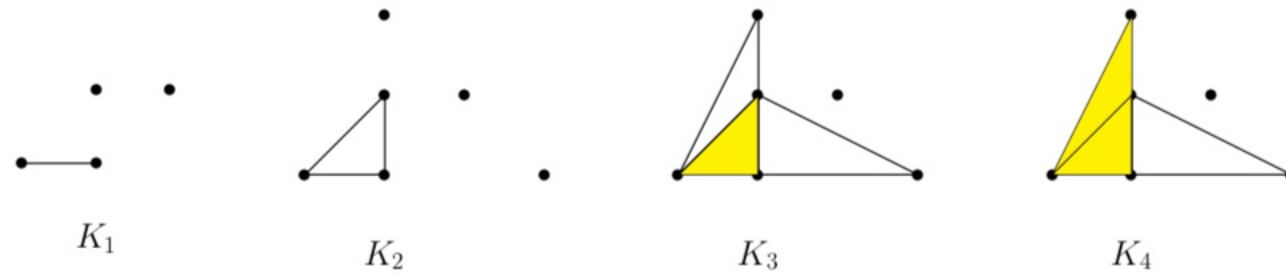




## Step 2: Point cloud to nested complexes



# Step 3: Nested complexes to barcode



$$\beta_0(K_1) = 3$$

$$\beta_1(K_1) = 0$$

$$\beta_0(K_2) = 4$$

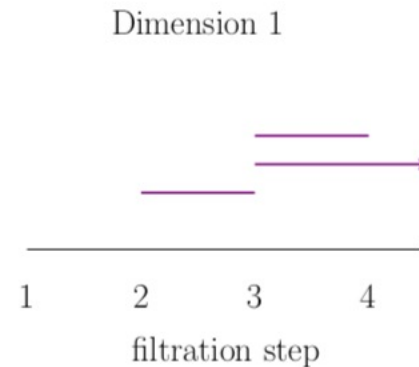
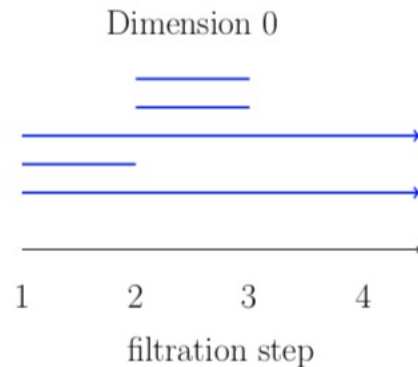
$$\beta_1(K_2) = 1$$

$$\beta_0(K_3) = 2$$

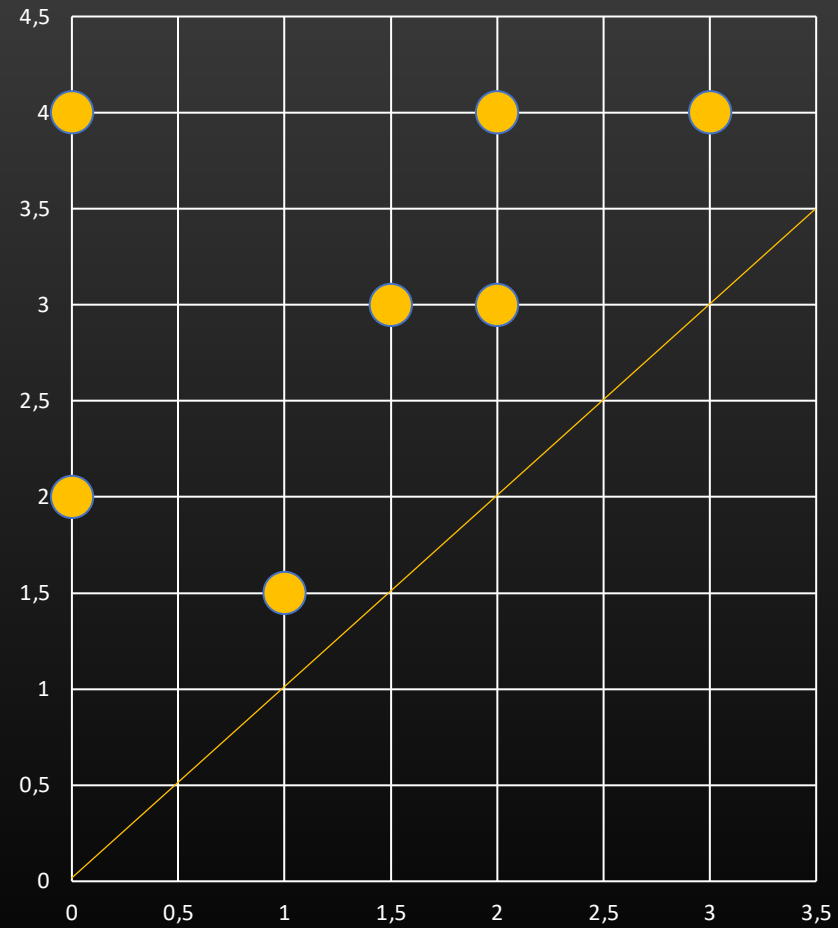
$$\beta_1(K_3) = 2$$

$$\beta_0(K_4) = 2$$

$$\beta_1(K_4) = 1$$



# Barcodes vs persistence diagrams (PD)



# Stability

- The set of barcodes/persistence diagrams can be equipped with a variety of **earthmover-type distances**: the Wasserstein distances of  $L_p$ -type and the bottleneck distance of  $L_\infty$ -type.
- Most reasonable known instantiations of the TDA pipeline are **Lipschitz continuous** with respect to Hausdorff distance on point clouds and bottleneck distance on persistence diagrams.

# Practicalities

- There are extensive **libraries of software**, mostly open source, for TDA computations (e.g., GUDHI, Ripser, Flagser, **Giotto-TDA**,...).
- There exist **“inverse analysis” tools** for interpreting results of TDA computations (e.g., work of Hiraoka et al.).

# From one to many parameters

- In real data, there are often several parameters along which it would be natural to filter (e.g., some notion of density or time).
- Generalization from one to many parameters poses serious problems, for reasons of both theory and implementation: in general, there is **no analogue of barcodes or persistence diagrams**.
- Common approaches for two parameters
  - Restrict to lines in the plane determined by the two parameters: **fibered bar code**.
  - Focus on **decompositions into blocks** (instead of bars) when possible.

Static TDA input to ML

# Strategies for vectorization/featurization

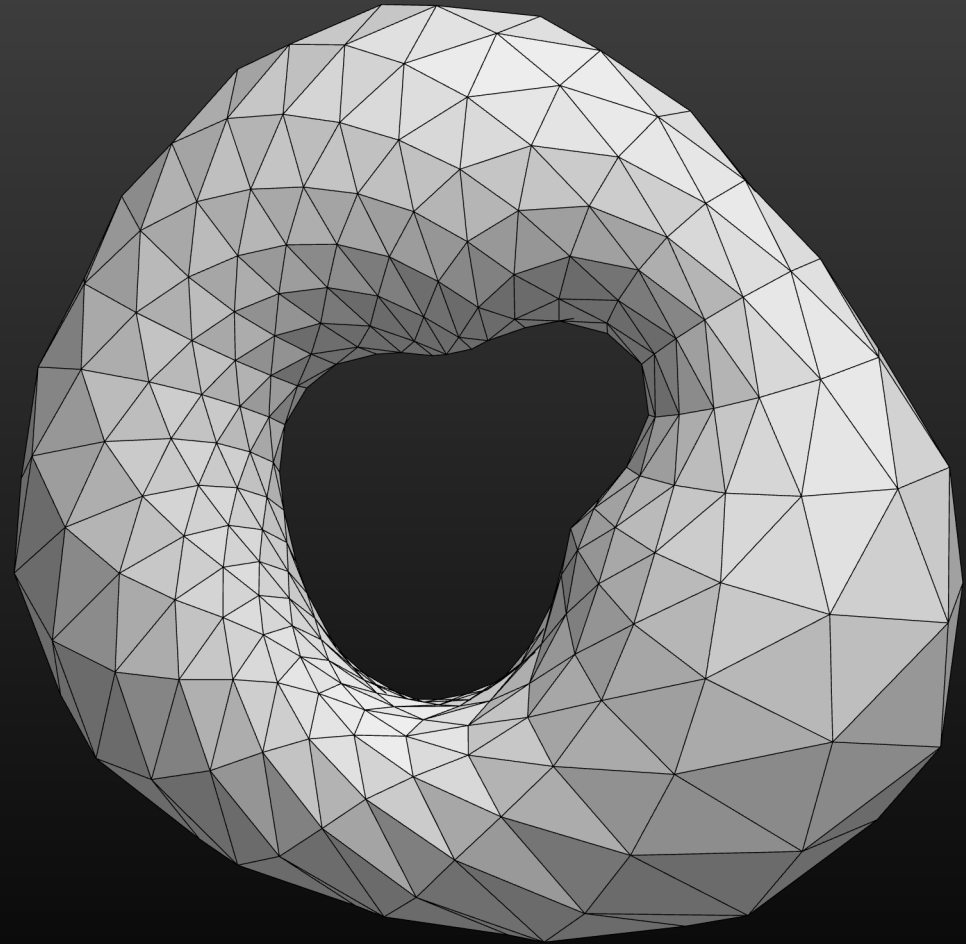
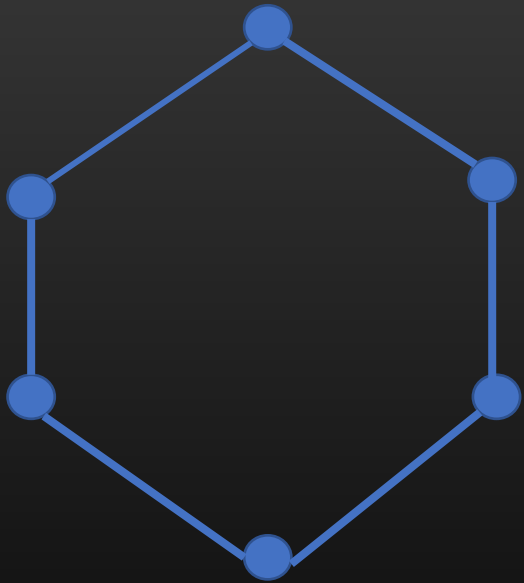
- **Problem:** Cannot compute statistics in the space of barcodes or the space of persistence diagrams.
- **Solution:**
  - Define a Lipschitz-continuous mapping from the space of barcodes/persistence diagrams to a vector space  $\mathcal{V}$  equipped with an inner product.
  - Compute statistics in  $\mathcal{V}$ !
  - Two main types:
    - Embeddings into finite-dimensional Euclidean spaces
    - Kernel methods: defining generalized scalar product on PD, i.e., see PD as elements of a Hilbert space

Few trainable parameters

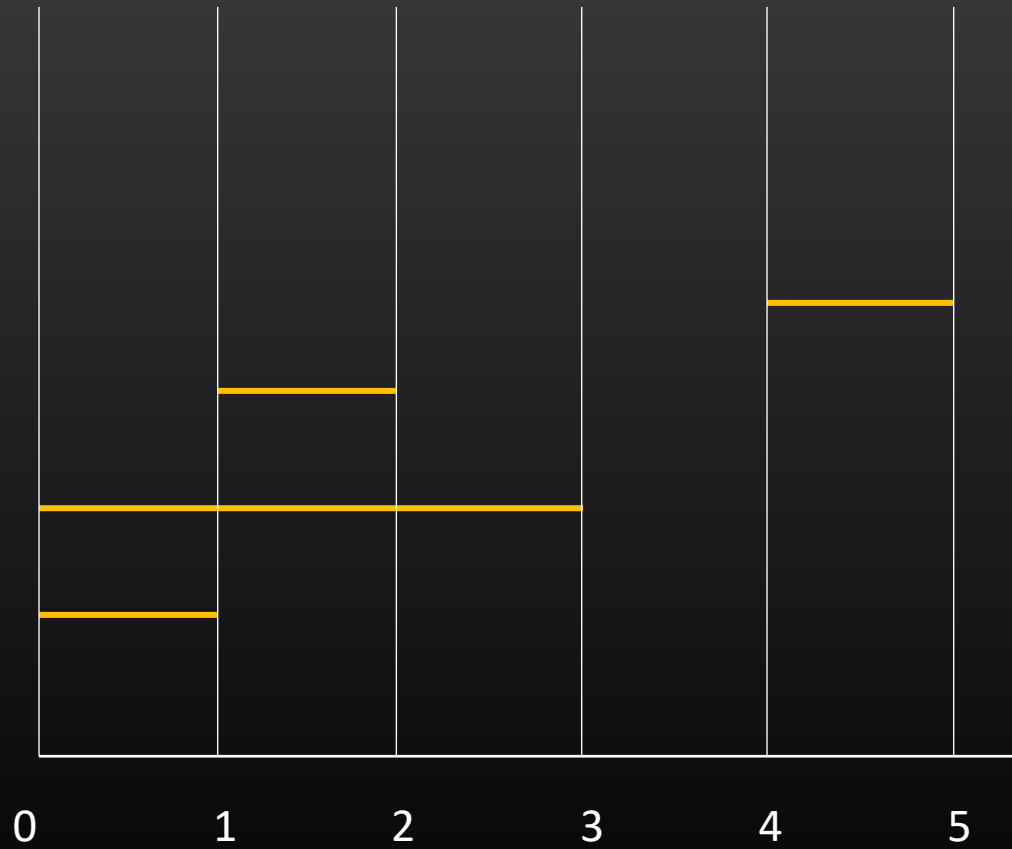
Expensive



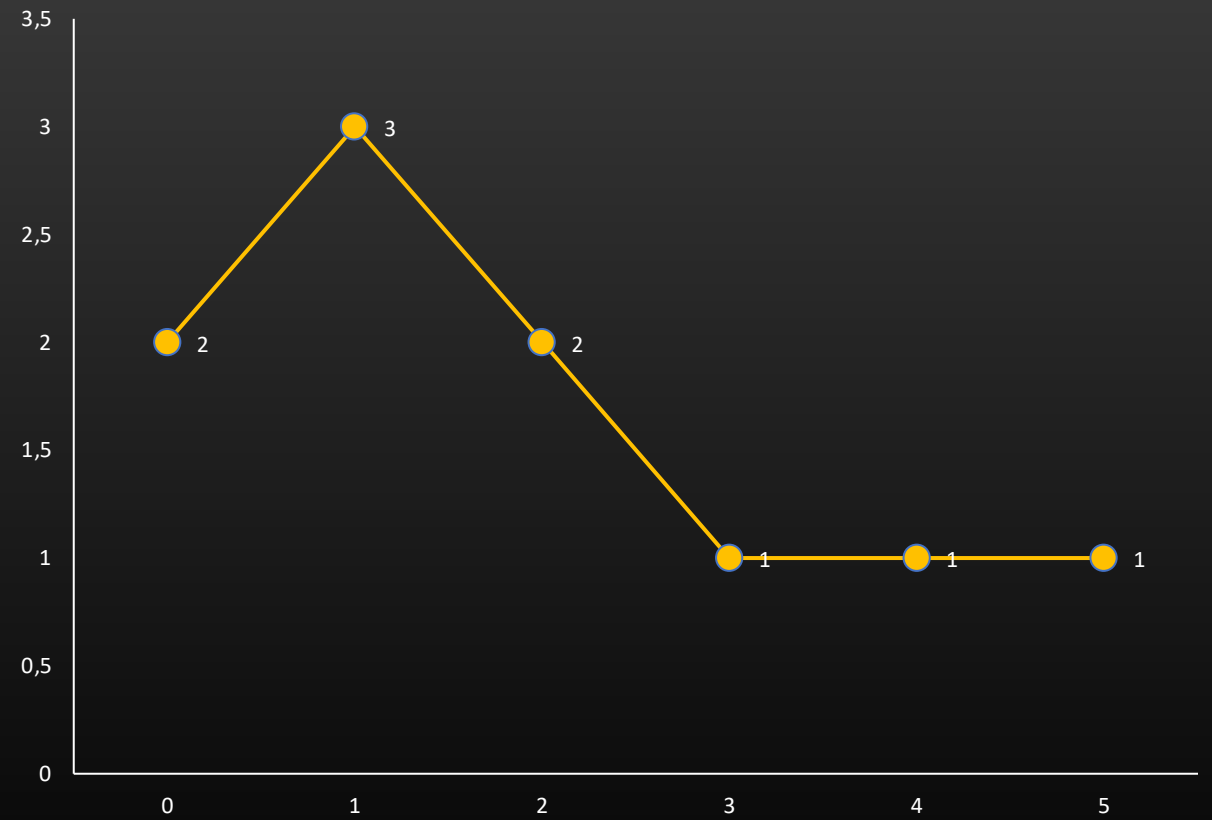
# Cavities



# Betti curves

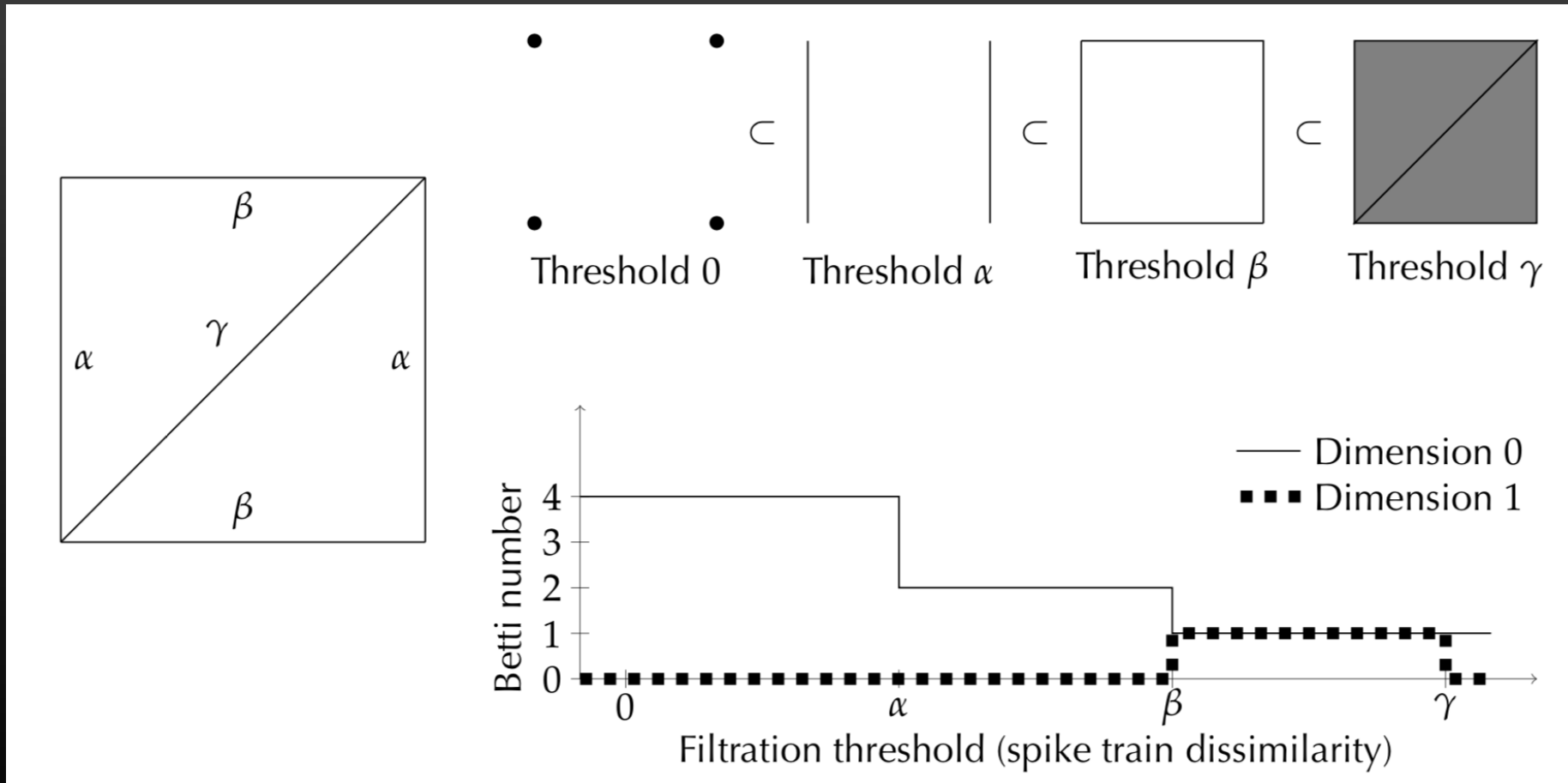


Bar code for cavities of dimension  $k$

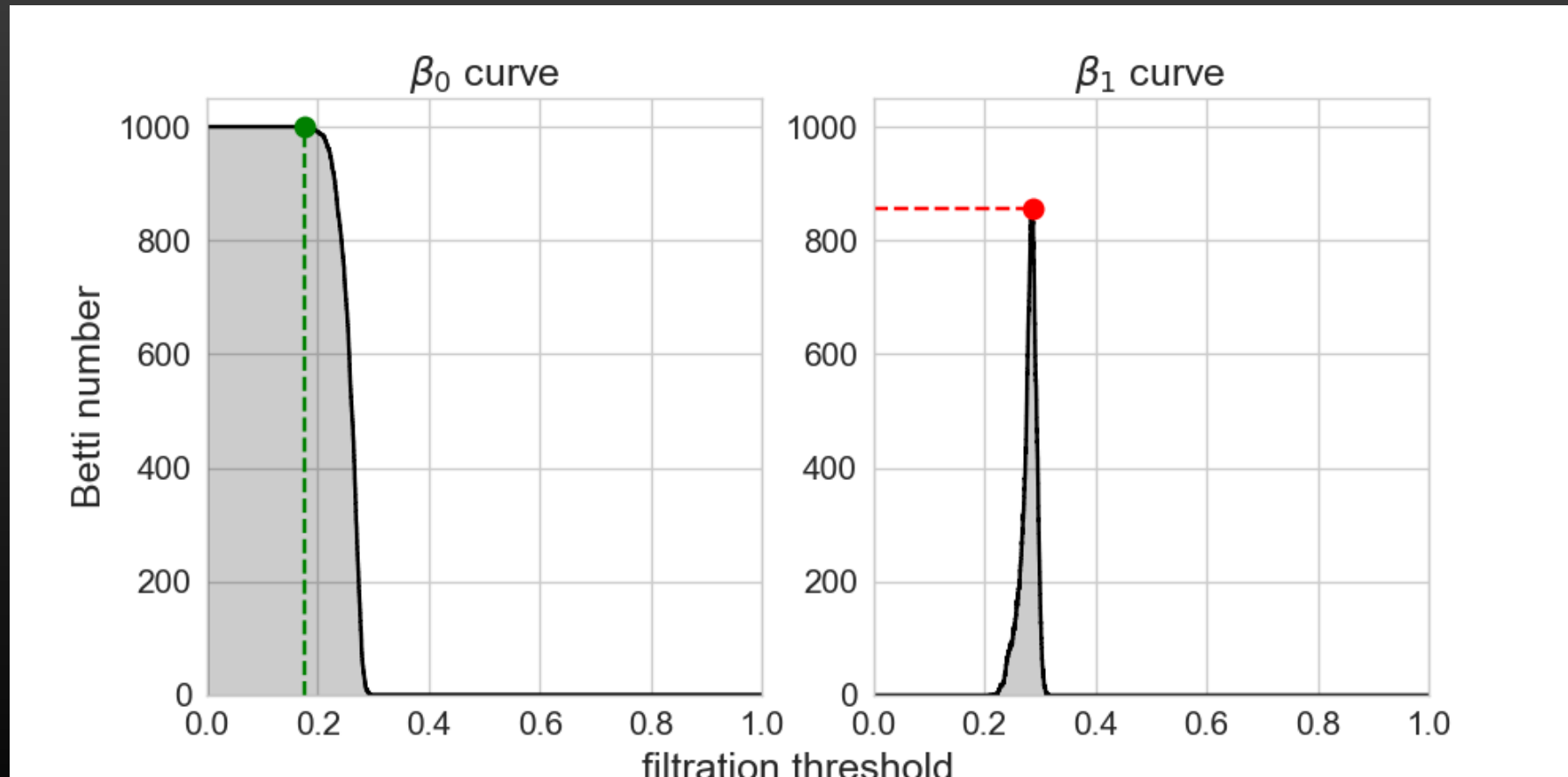


Betti\_k curve

# Nested complex to Betti curve

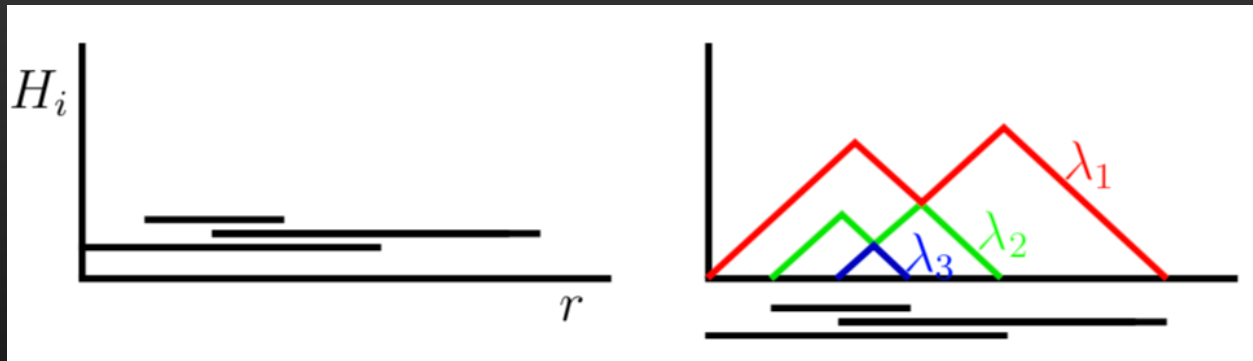


# Extracting numerical features



# Persistence landscapes

- Barcodes also give rise to **persistence landscapes**.



$$\lambda = \{ \lambda_k : \mathbb{R} \rightarrow \mathbb{R} \cup \{ \infty \} \mid k \in \mathbb{N} \}$$

- The **L2-landscape distance** between barcodes  $B$  and  $B'$  with associated landscapes  $\lambda$  and  $\lambda'$ :

$$\Lambda(B, B') = \|\lambda - \lambda'\|_2 = \sum_{k=1}^{\infty} \left( \int |\lambda_k(t) - \lambda'_k(t)|^2 dt \right)^{\frac{1}{2}}$$

# Persistence curves

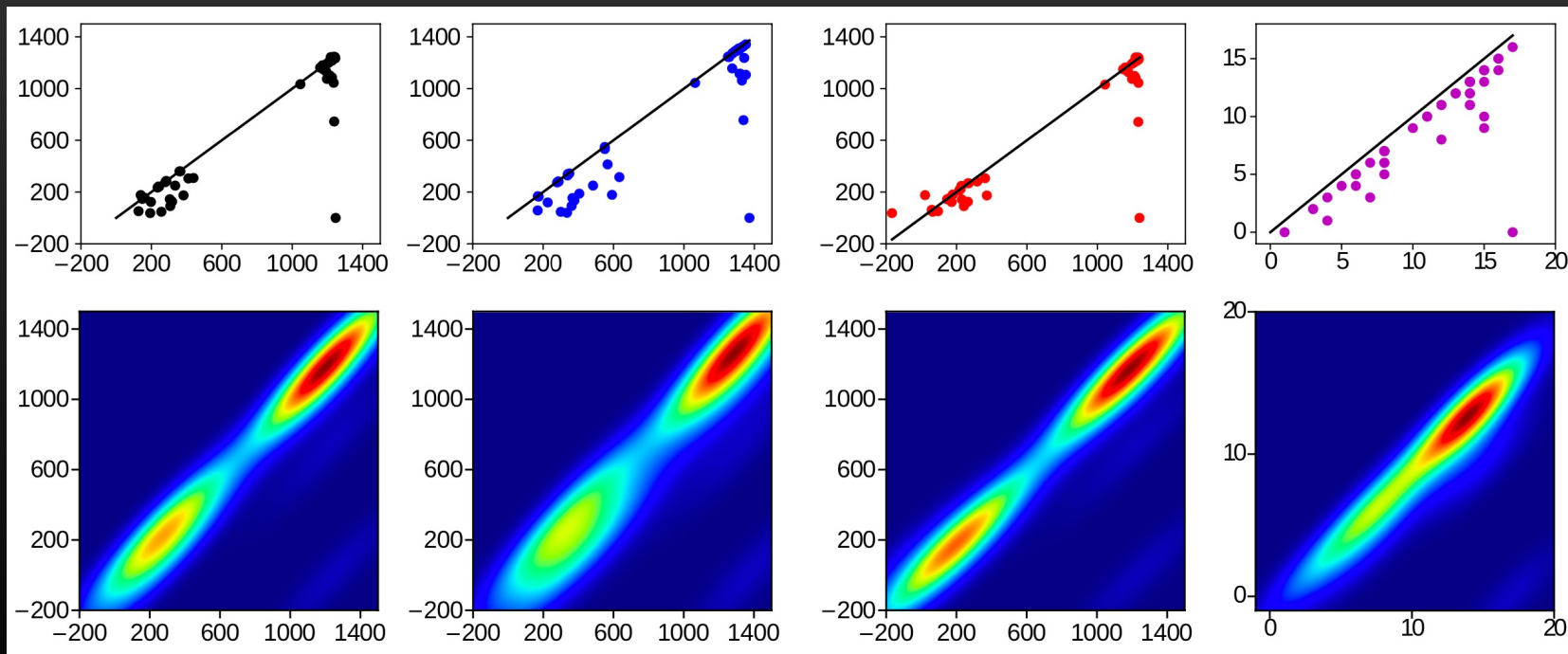
Name	Notation	$\psi(b, d, t)$	T
Betti	$\beta(D)$	1	sum
Midlife	$\mathbf{ml}(D)$	$(b + d)/2$	sum
Life	$\ell(D)$	$d - b$	sum
Multiplicative Life	$\mathbf{mul}(D)$	$d/b$	sum
Life Entropy [2]	$\mathbf{le}(D)$	$-\frac{d-b}{\sum(d-b)} \log \frac{d-b}{\sum(d-b)}$	sum
Midlife Entropy	$\mathbf{mle}(D)$	$-\frac{d+b}{\sum(d+b)} \log \frac{d+b}{\sum(d+b)}$	sum
Mult. Life Entropy	$\mathbf{mule}(D)$	$-\frac{d/b}{\sum(d/b)} \log \frac{d/b}{\sum(d/b)}$	sum
$k$ -th Landscape [5]	$\lambda_k(D)$	$\min\{t - b, d - t\}$	$\max_k$

Simultaneous **generalization of Betti curves and persistence landscapes**. Robust to input noise, efficient to compute, interpretable, and allowing weighting of relative importance of different regions in the PD.

For each  $t$ , compute  $\mathbb{T}(\{\psi(b, d, t) \mid b \leq t, d > t\})$ .

# Persistence images

- Smooth the PD: replace each point by a Gaussian kernel, then sum
- Discretize



(Image from Kanari, et al., Neuroinformatics, 2018.)

Adams et al., JMLR 2017

# ML methods applied to vectorized TDA

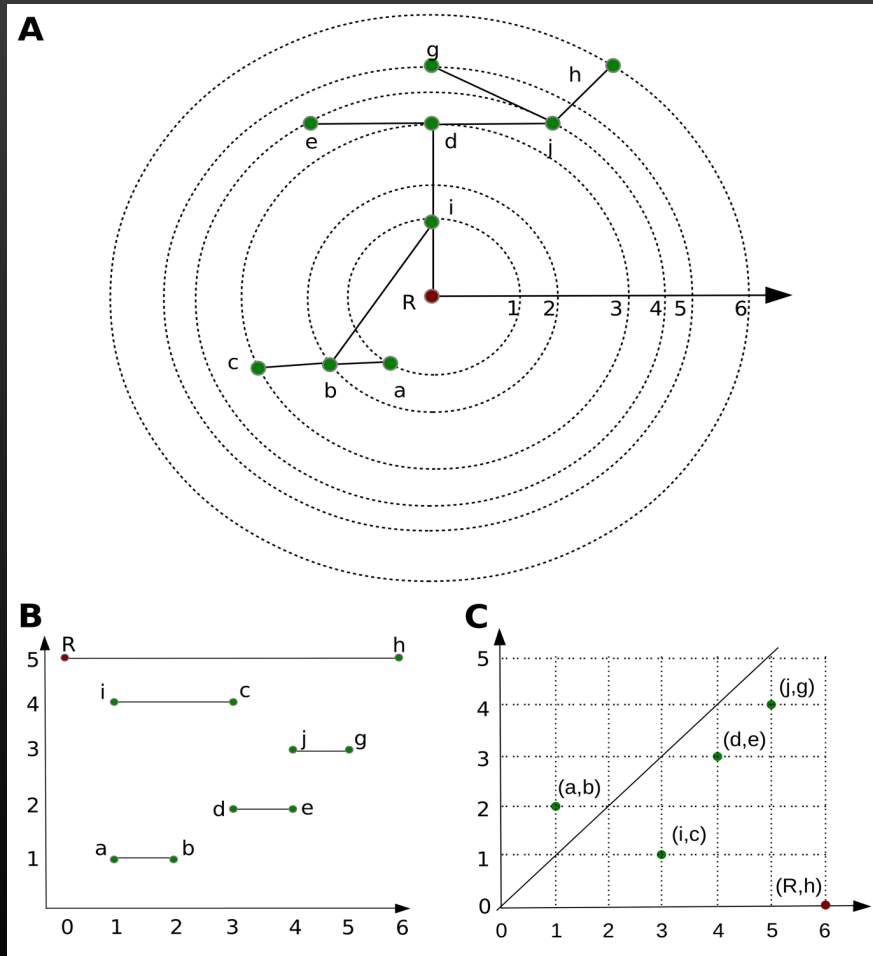
- Decision tree
- Random forest
- Support Vector Machine
- CNN
- GNN

Also possible to integrate a TDA layer into an ML model!



Applications

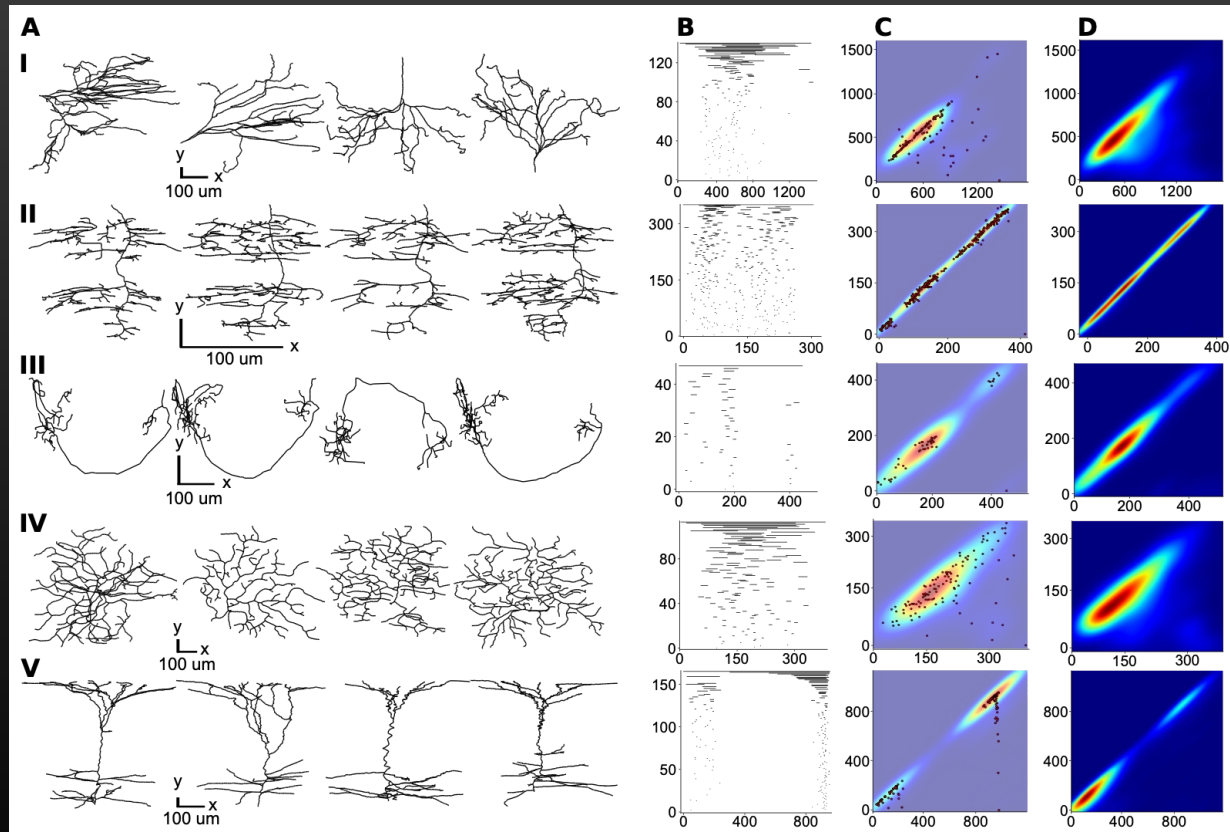
# Classification of neuron morphologies



**Idea:** Starting at the leaves and descending recursively to the root, decompose the tree into branches, while respecting the **Elder Rule**, i.e., at any bifurcation, the elder (longer) branch survives and the younger branch is broken off.

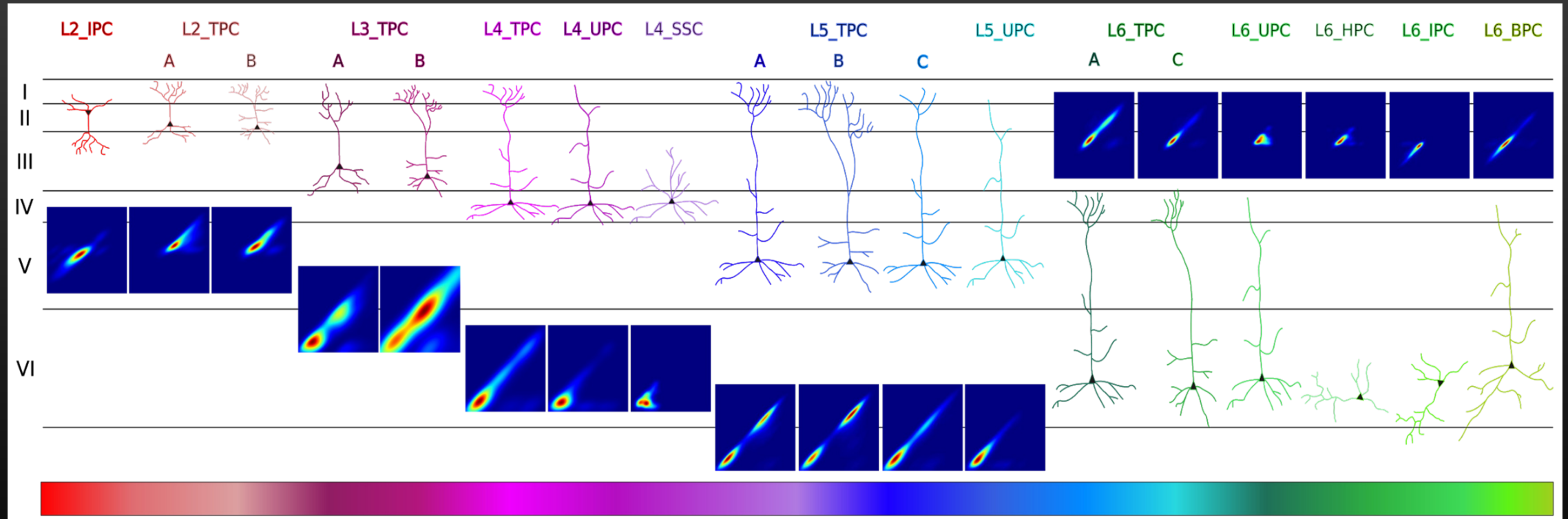
Integrate the **topology of the tree** and the **geometry of its embedding** in space into a surprisingly powerful **global descriptor**.

# Classification of neuron morphologies

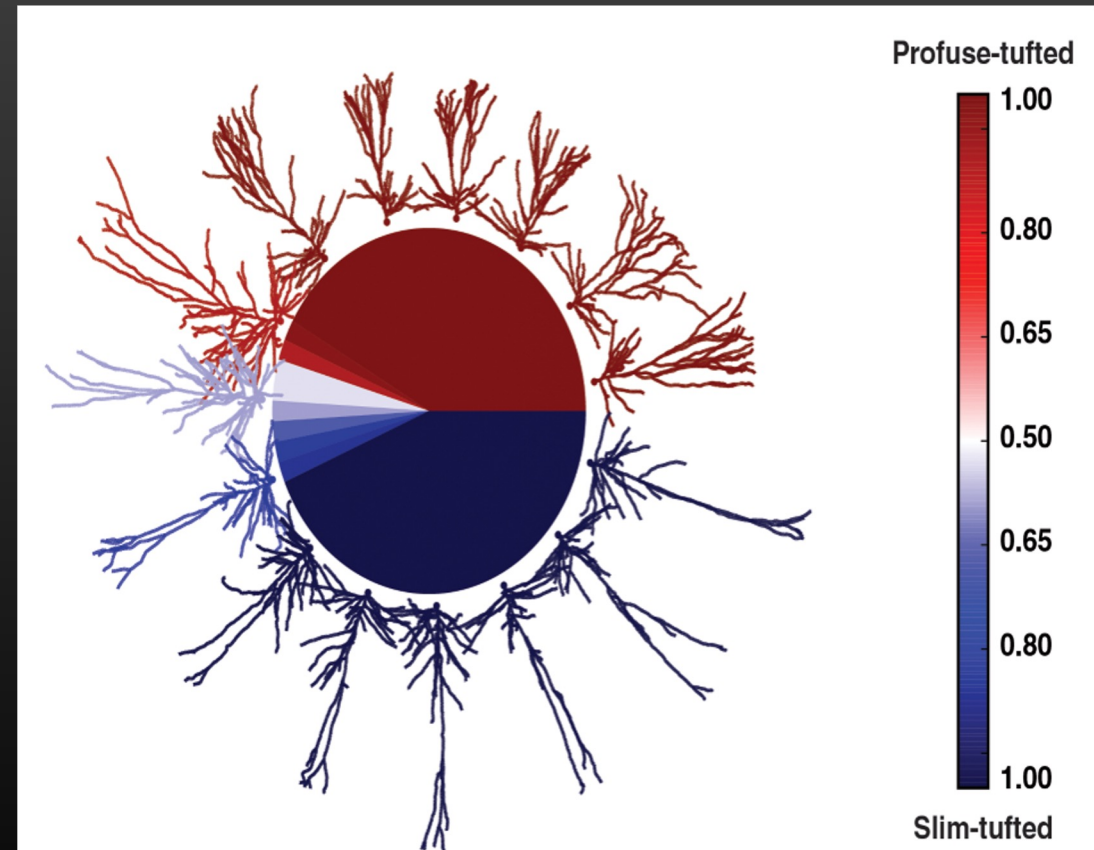
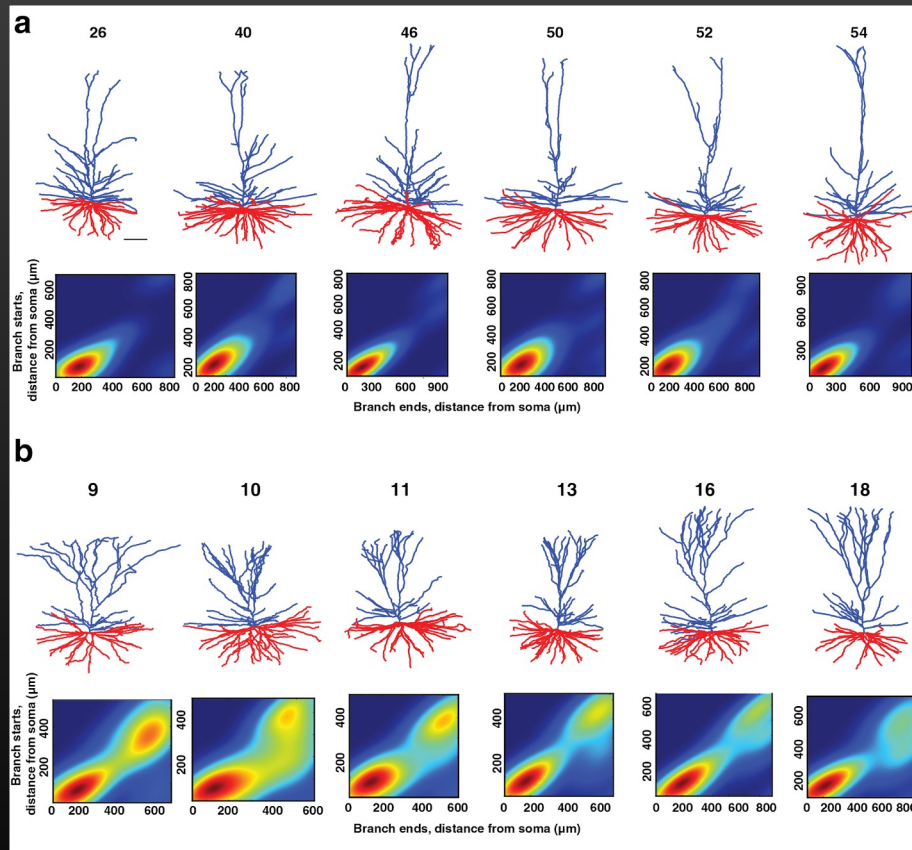


Kanari et al., Neuroinformatics 2017  
Kanari et al., Cerebral Cortex 2019

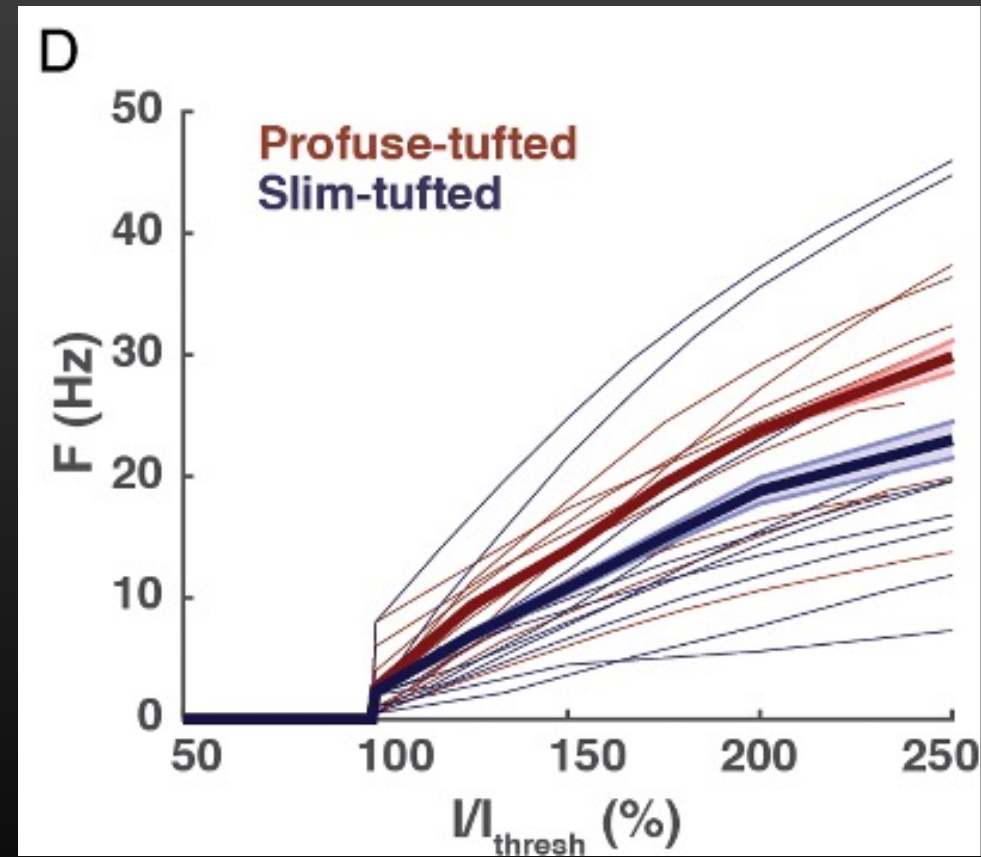
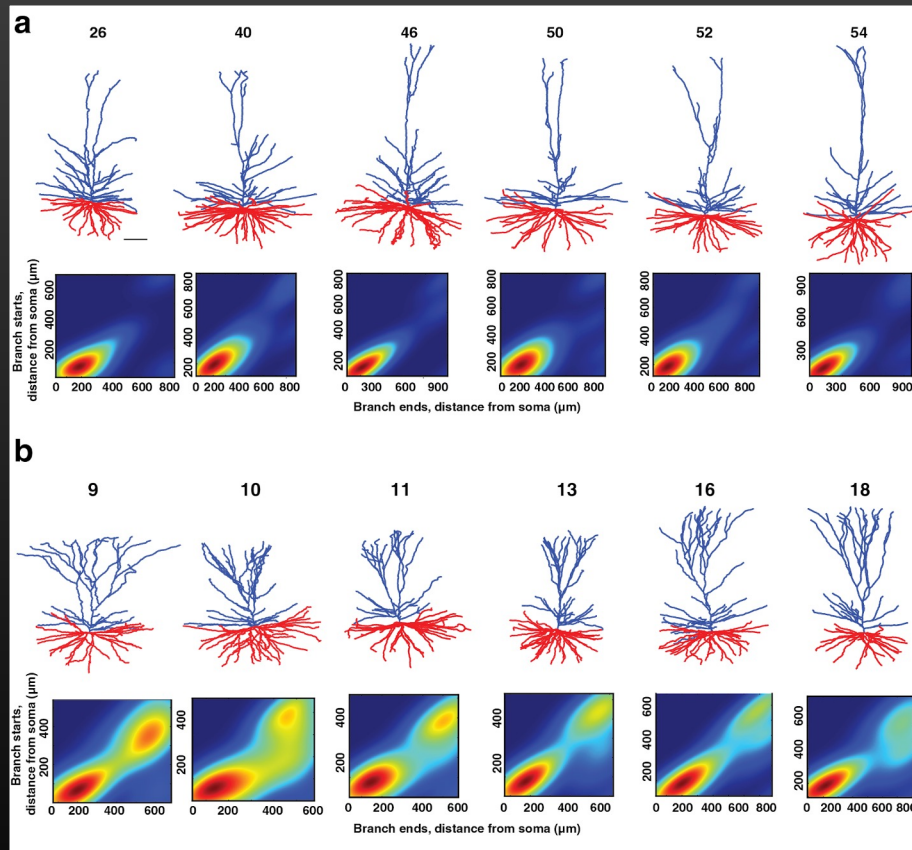
# Classification of neuron morphologies



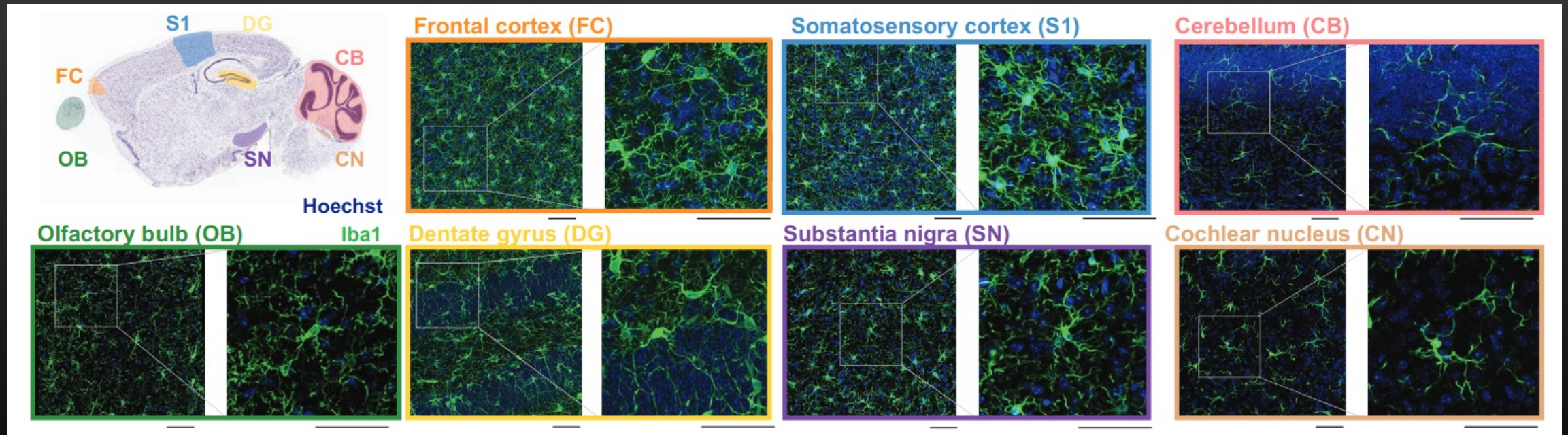
# Classification of neuron morphologies



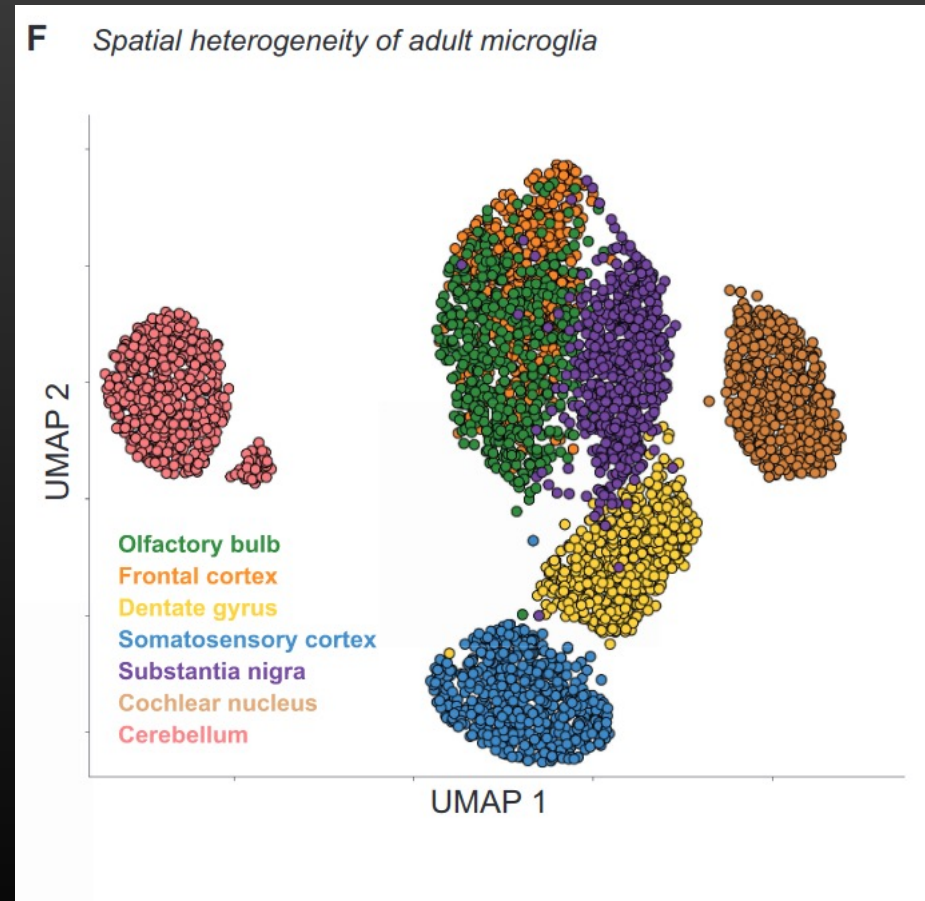
# Classification of neuron morphologies



# Classification of microglia

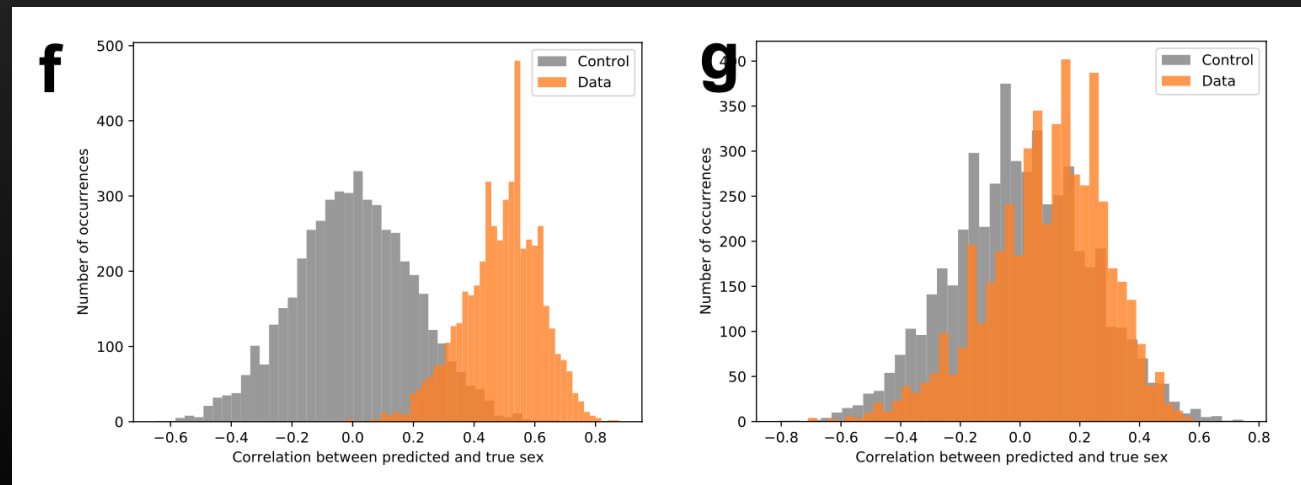
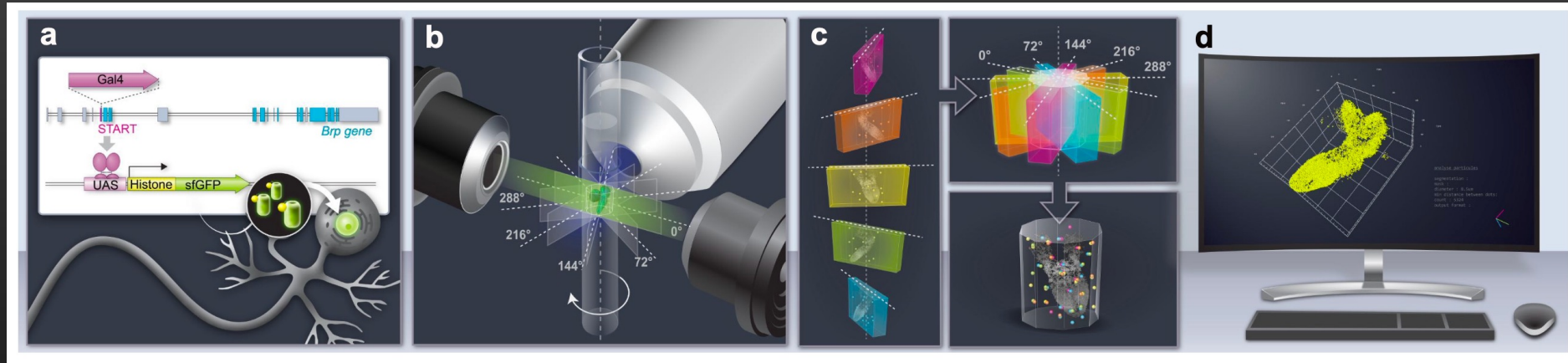


# Classification of microglia



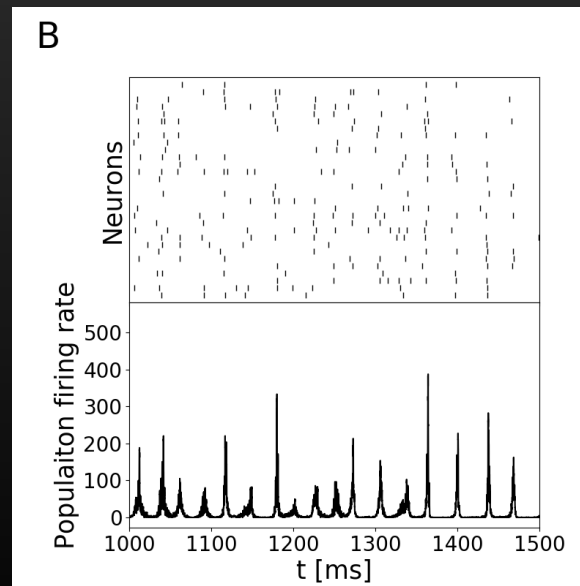
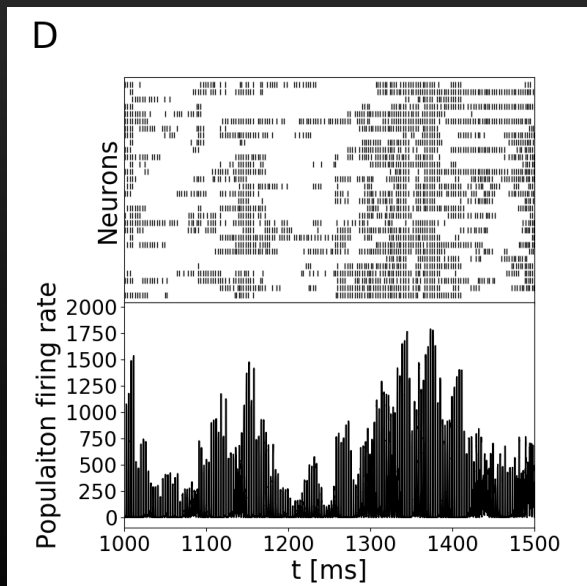


# Sexual dichotomy in larval fruitflies



# Classification of neural dynamics

- For a range of activity parameters, associate to an active Brunel network a **weighted graph**, to which we apply tools of **persistent homology**.
- Extract simple topological features of each dynamic regime.
- Use these to train a (highly accurate!) **classifier**.



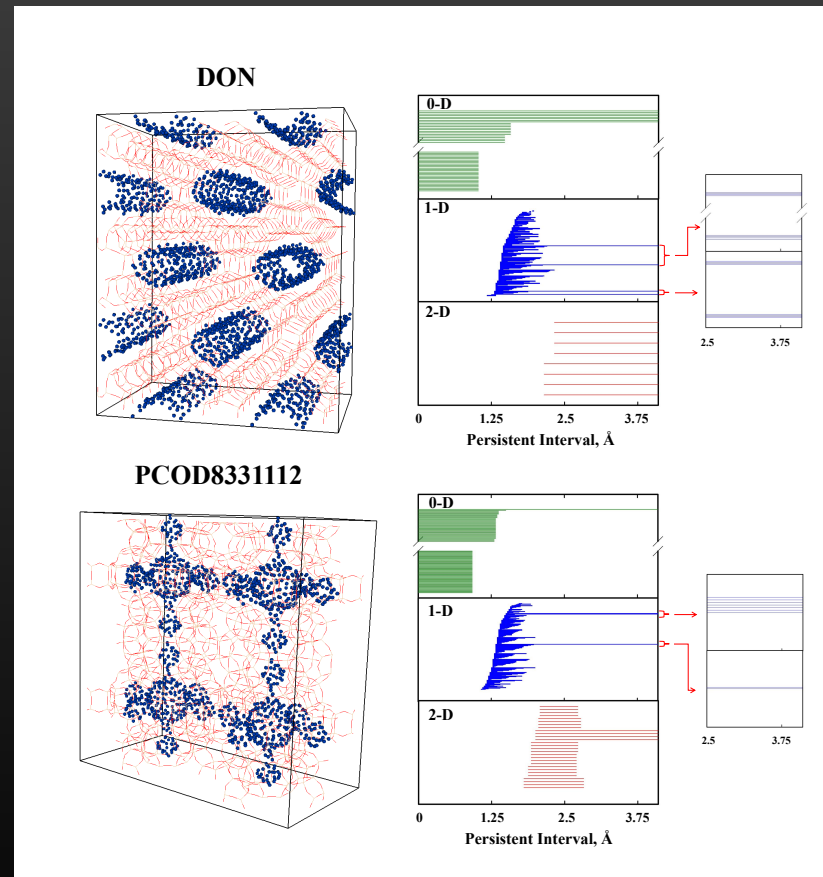
# Classification of neural dynamics

## Automated classification of network dynamics

- For a range of activity parameters, associate to an active Brunel network a **weighted graph**, to which we apply tools of **persistent homology**.
- Extract simple topological features of each dynamic regime.
- Use these to train a (highly accurate!) **classifier**.

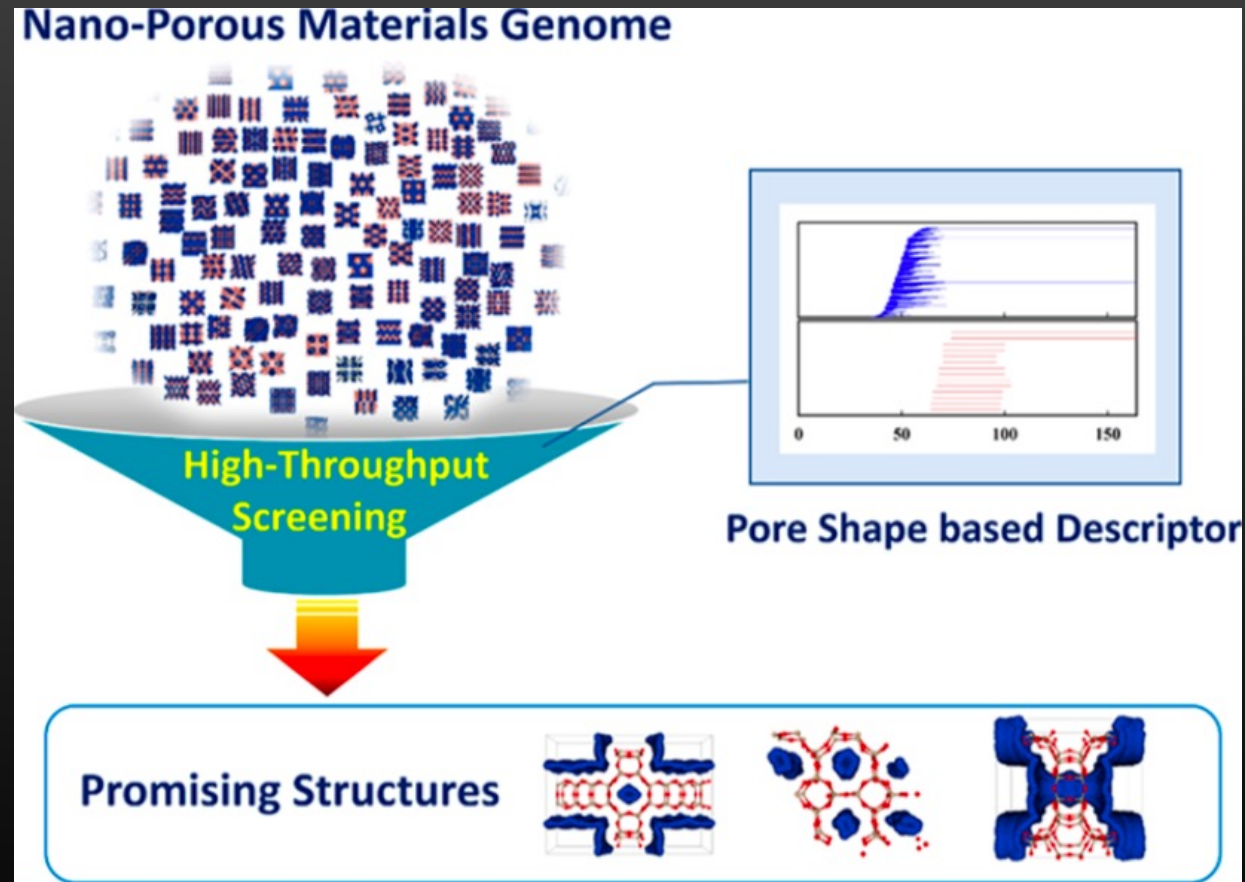
Training set	Testing set			
	Ver. 1	Ver. 2	Ver. 3	All ver.
Version 1	100% (28)	86.67% (180)	91.18% (170)	89.68% (378)
Version 2	97.69% (130)	100% (24)	93.33% (240)	95.18% (394)
Version 3	99.23% (130)	99.17% (240)	100% (24)	99.23% (394)
All versions	100% (28)	100% (24)	100% (24)	100% (76)

# Classification of nanoporous crystalline materials



Lee et al., Nature Communications 2017  
Lee et al., J Chem Thy Comput 2018

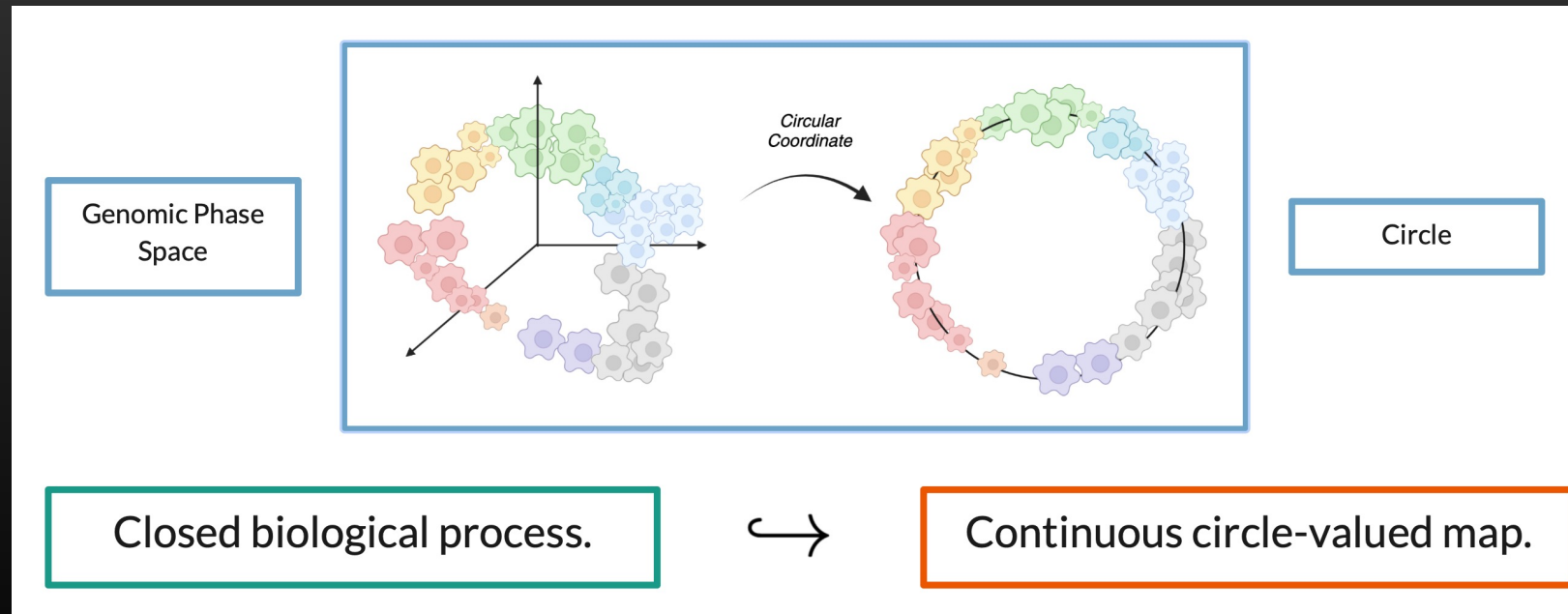
# Classification of nanoporous crystalline materials



Lee et al., Nature Communications 2017  
Lee et al., J Chem Thy Comput 2018

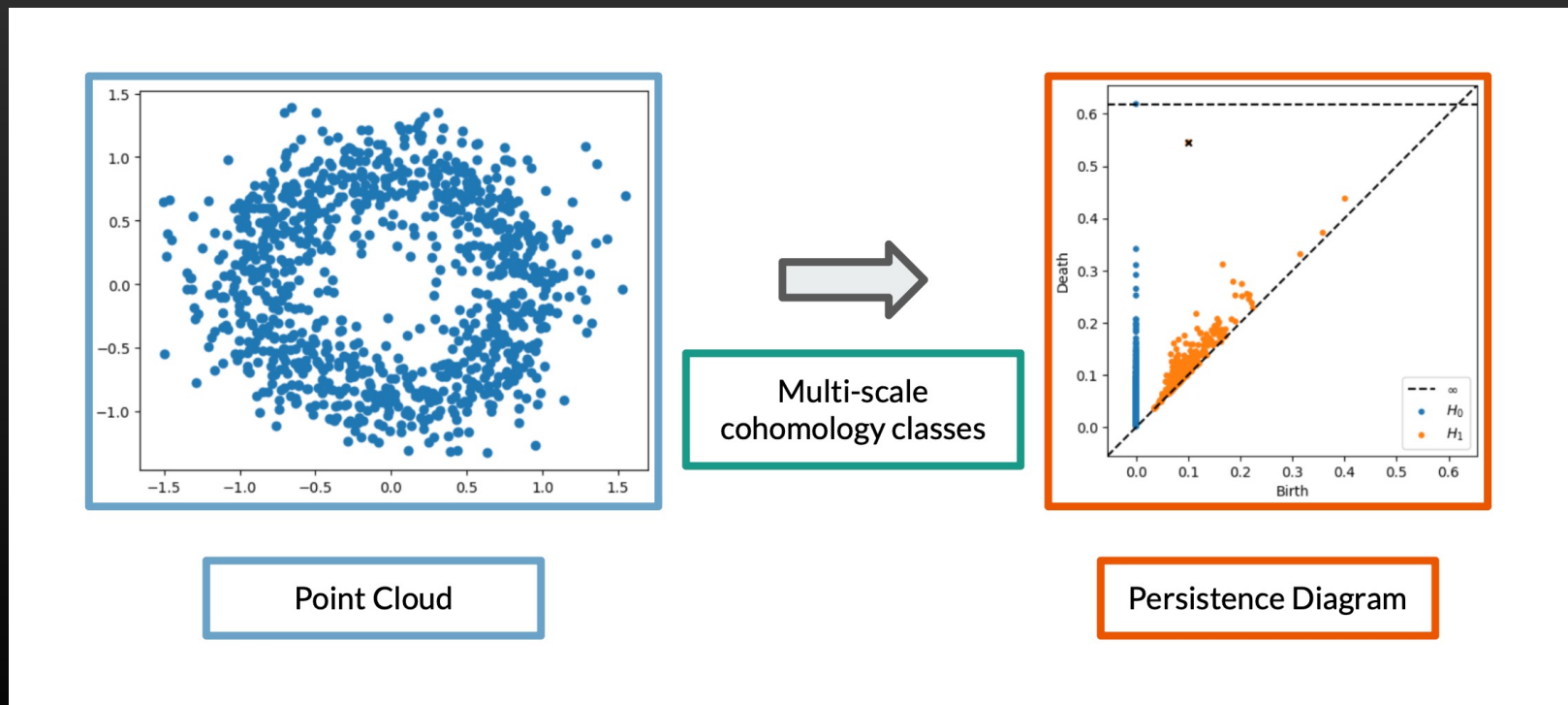
# Detection of gene cascades in single-cell data

- **Idea:** Generate biological hypotheses about closed processes in single-cell RNA seq data using topology and geometry



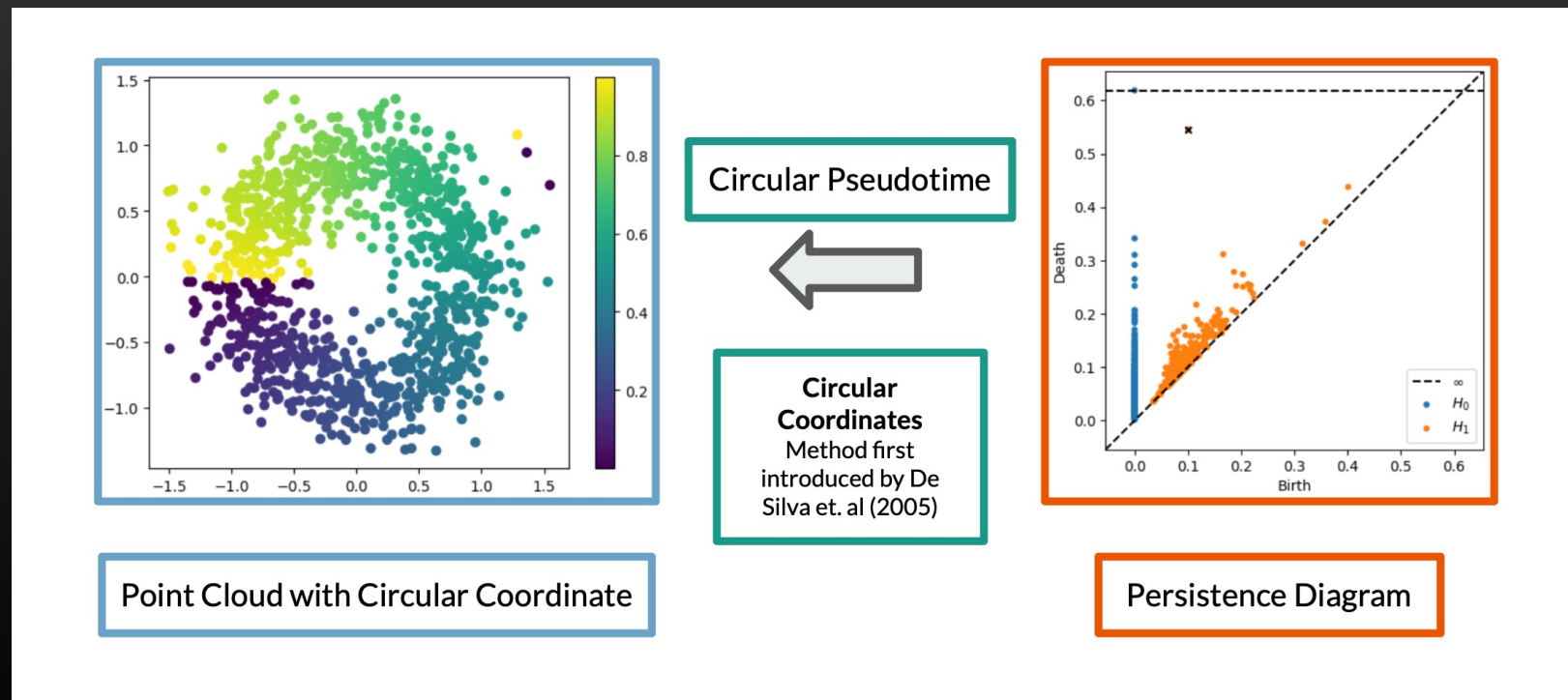
# Detection of gene cascades in single-cell data

- **Idea:** Generate biological hypotheses about closed processes in single-cell RNA seq data using topology and geometry



# Detection of gene cascades in single-cell data

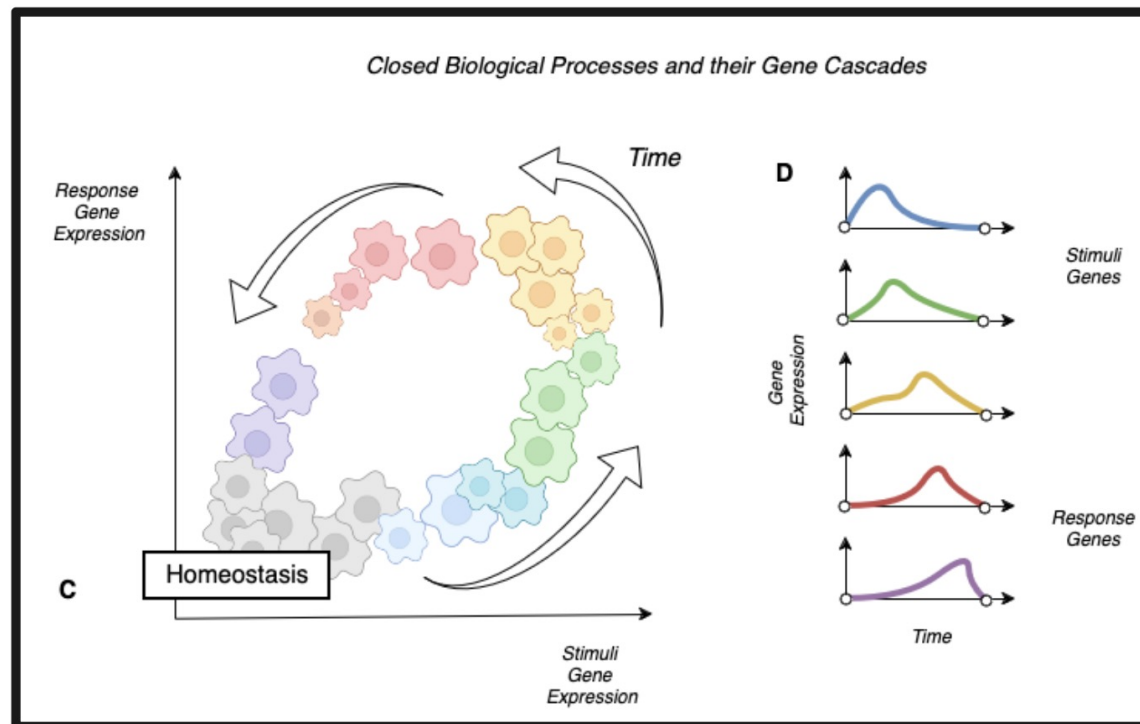
- **Idea:** Generate biological hypotheses about closed processes in single-cell RNA seq data using topology and geometry





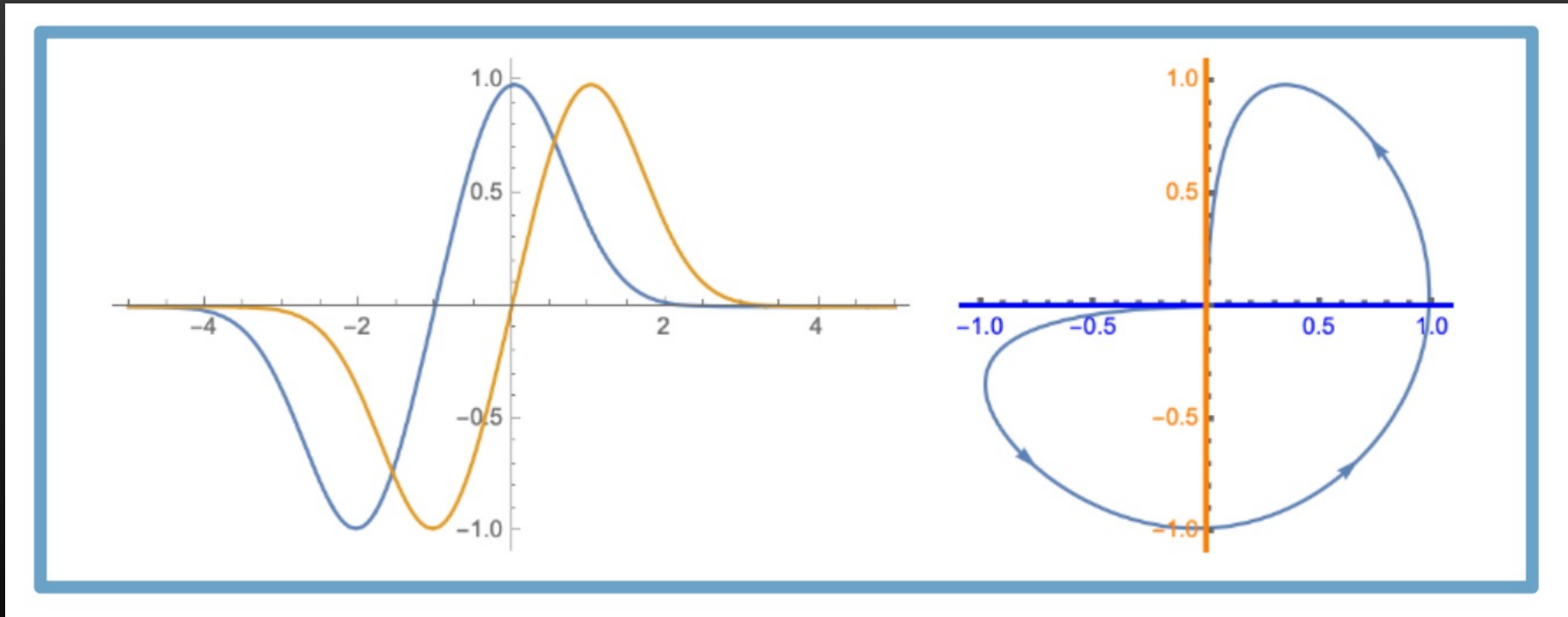
# Detection of gene cascades in single-cell data

Closed biological process

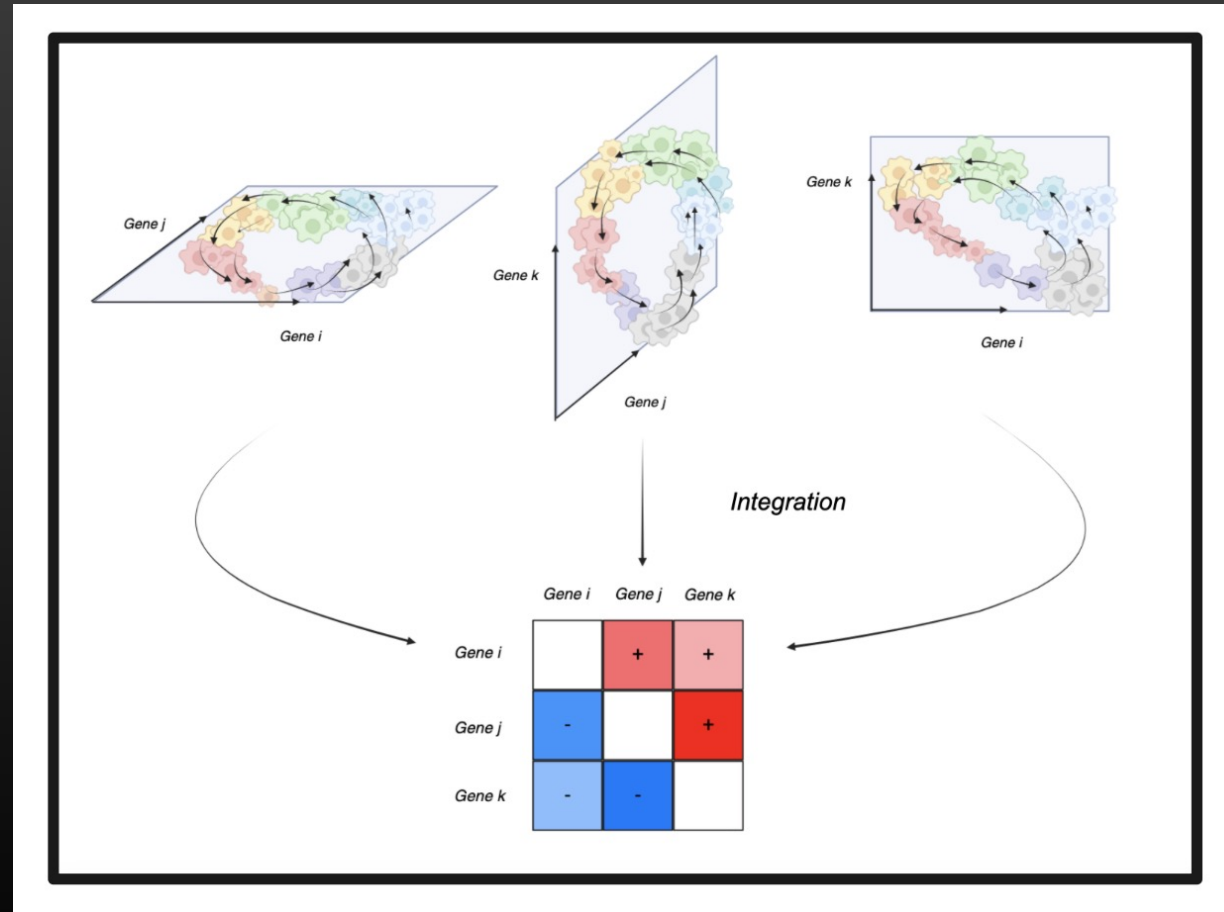


Gene expression cascade

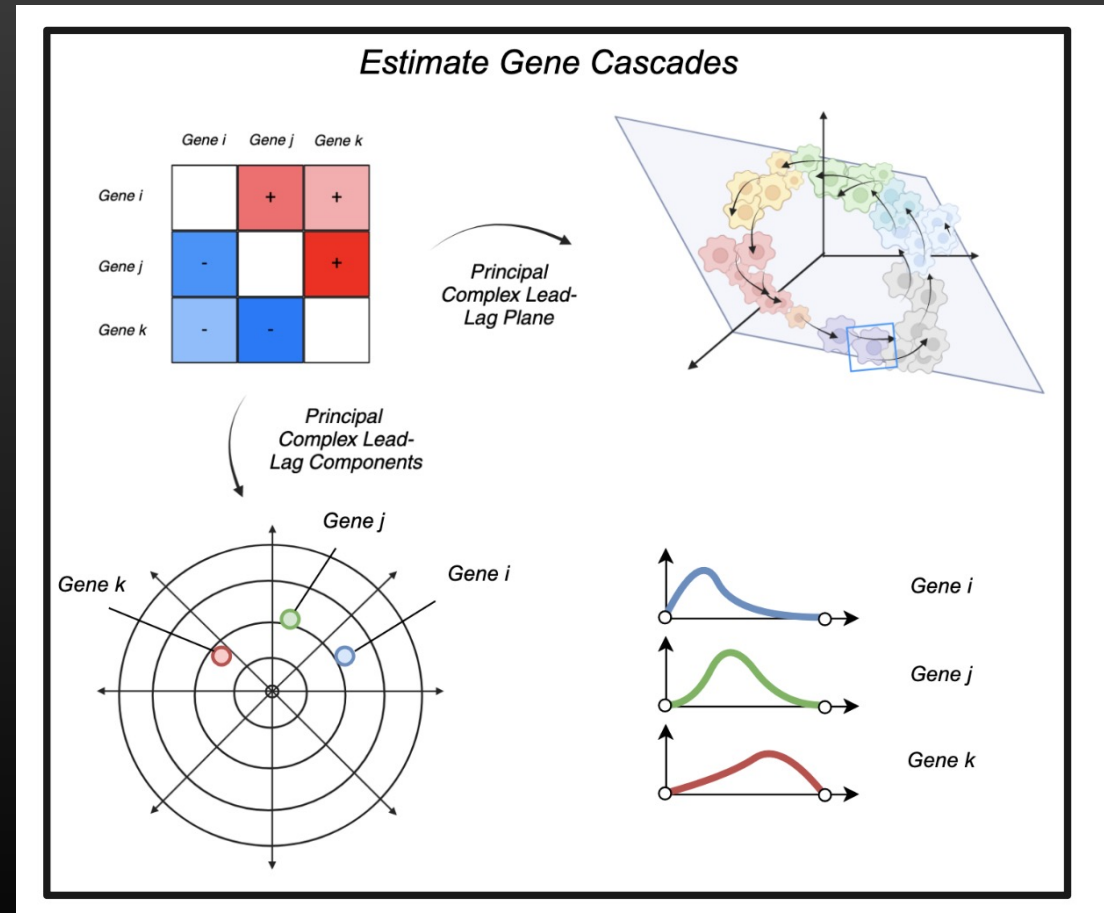
# Detection of gene cascades in single-cell data



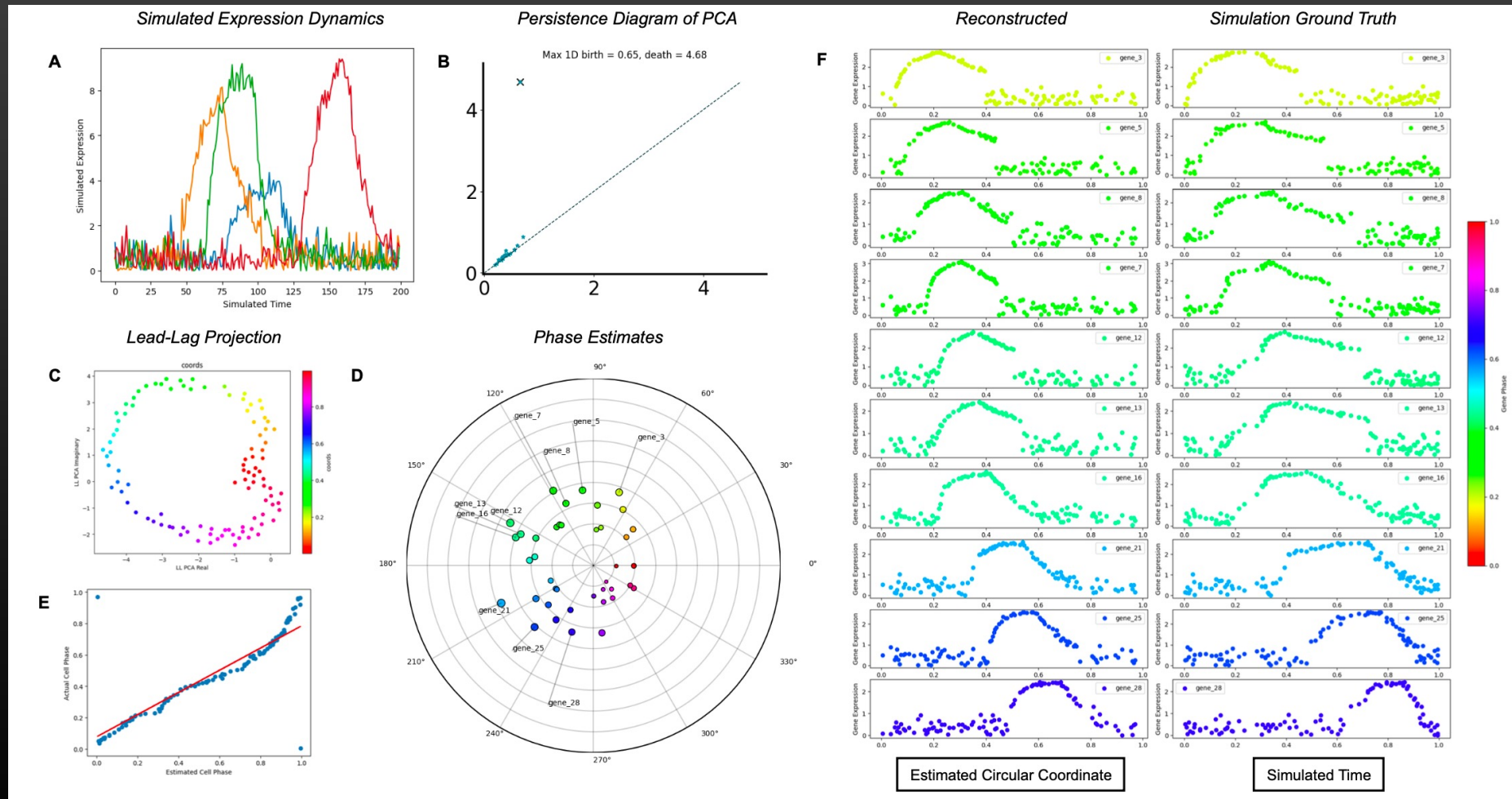
# Detection of gene cascades in single-cell data



# Detection of gene cascades in single-cell data



# Detection of gene cascades in single-cell data



Thank you to our funding agencies

- SNSF
- Innosuisse
- Blue Brain Project

*Merci !*



Members of the Laboratoire de  
Topologie et Neurosciences