

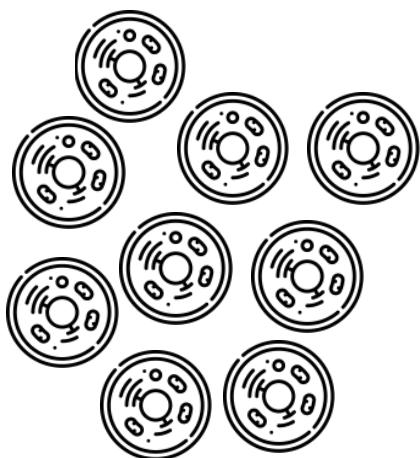
Clustering & cell annotation

Lieke Mchielsen

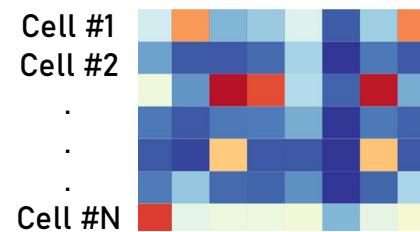
Department of Human Genetics, LUMC
Delft Biinformatics Lab, TU Delft

How can we identify cell populations?

Mystery cells

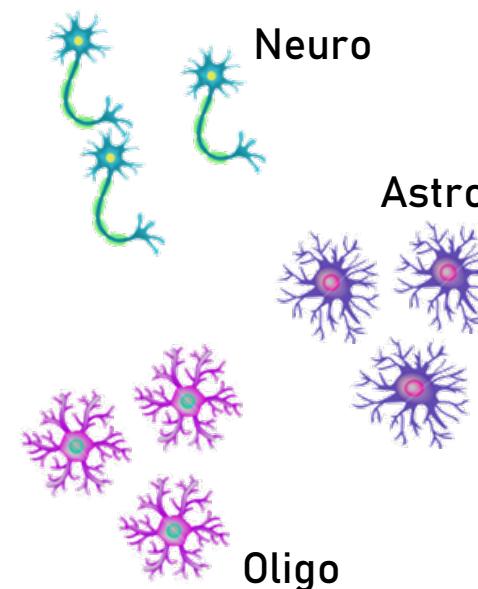


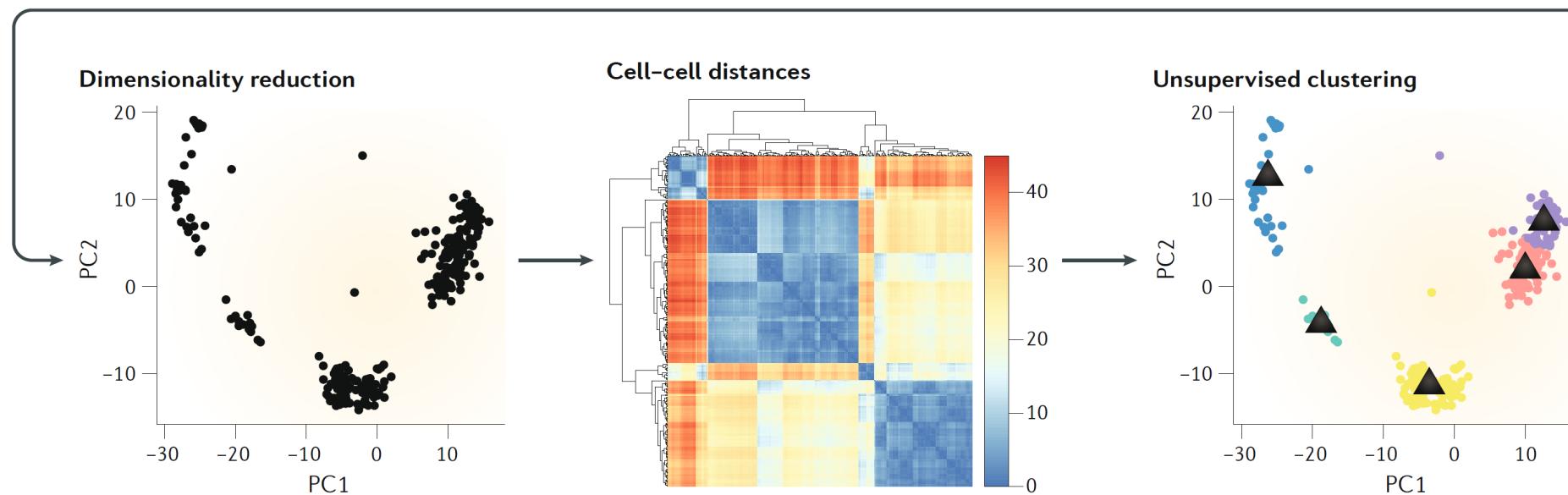
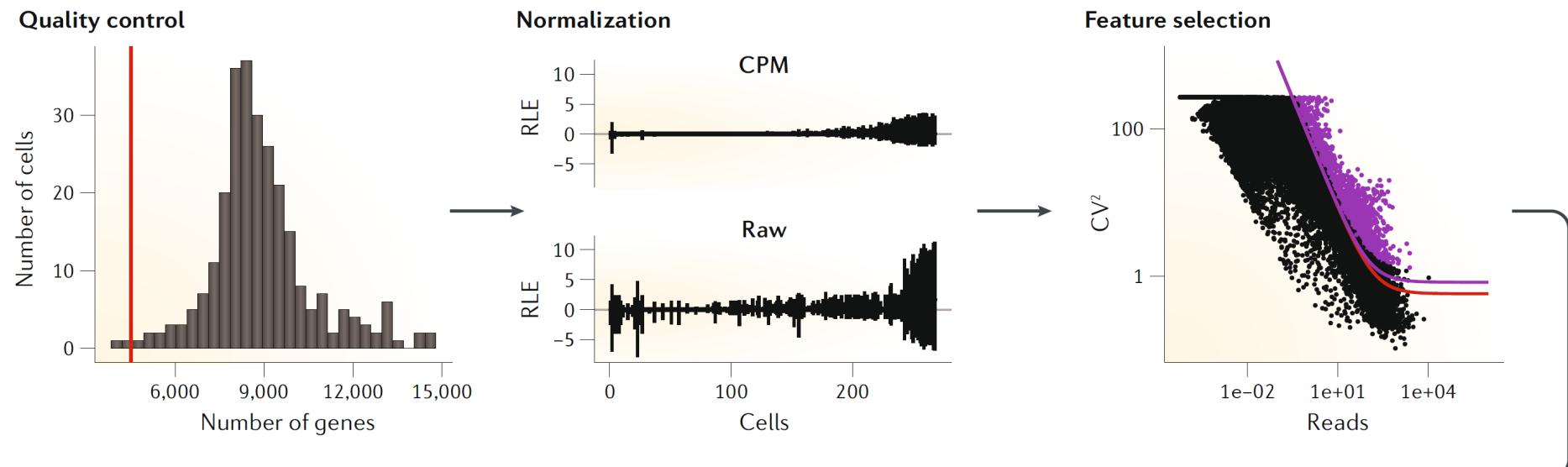
Measure

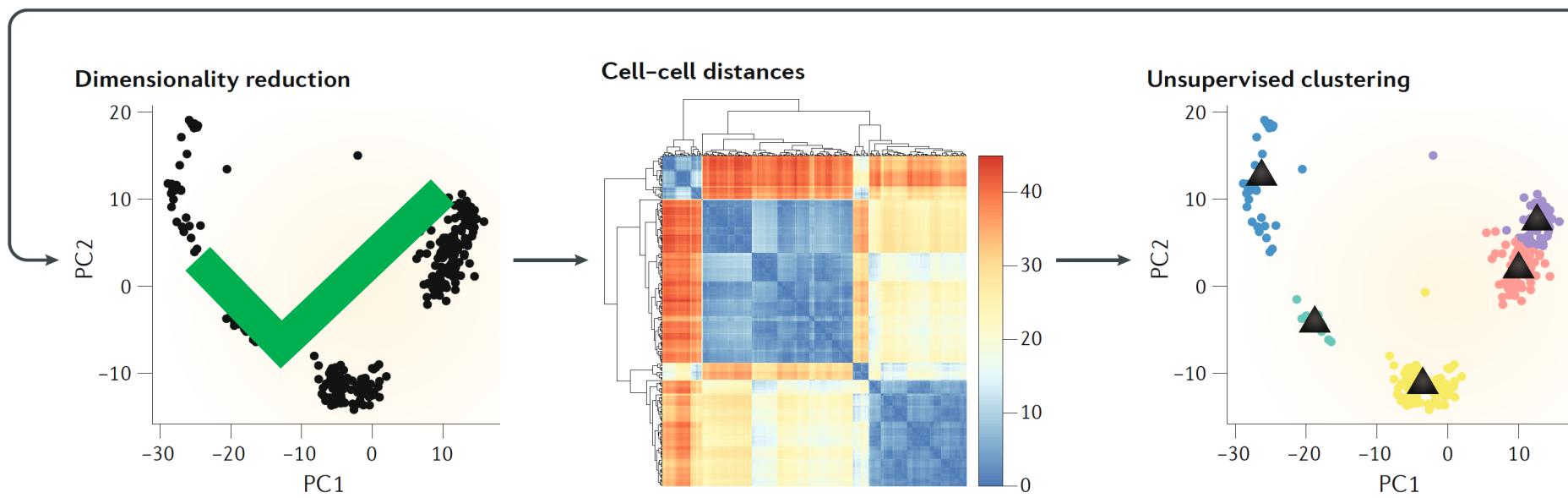
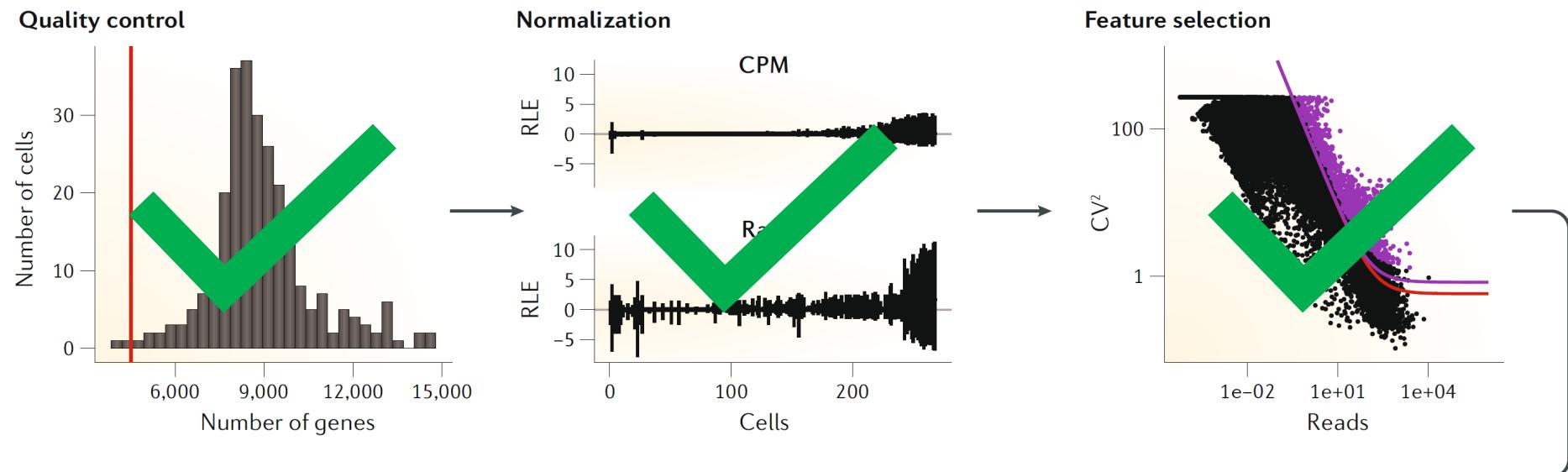


Identify

Cell populations







Outline

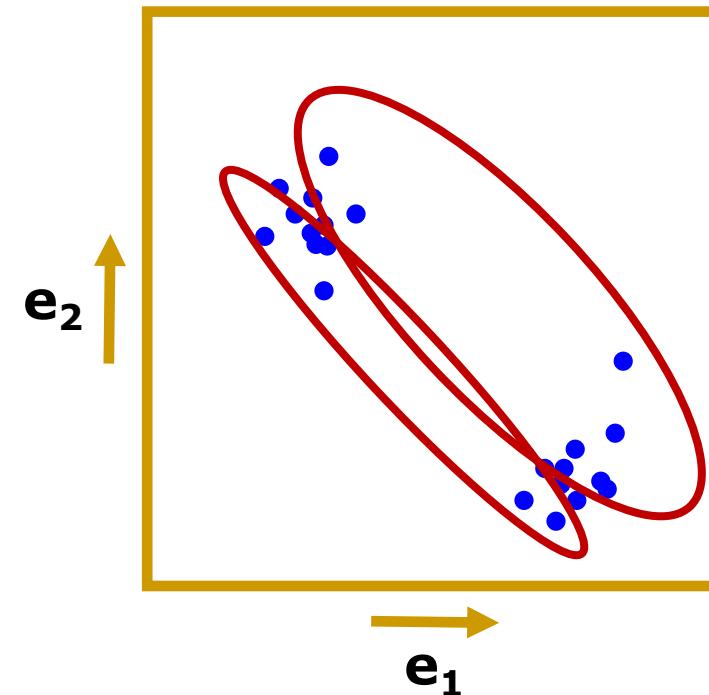
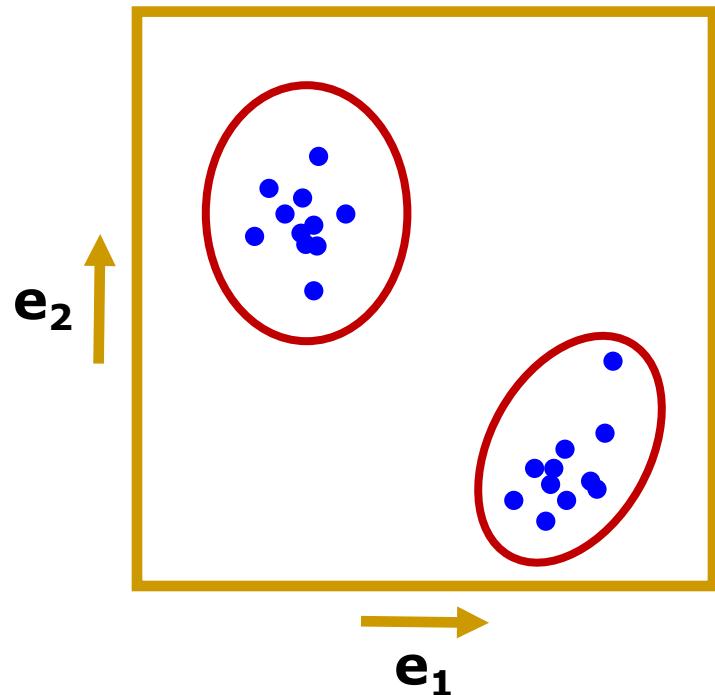
- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- scRNA-seq clustering
- Annotating clusters
- Cluster validation
- Challenges & outlook

Outline

- **Introduction to clustering**
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- scRNA-seq clustering
- Annotating clusters
- Cluster validation
- Challenges & outlook

Clustering

- What defines a good clustering?



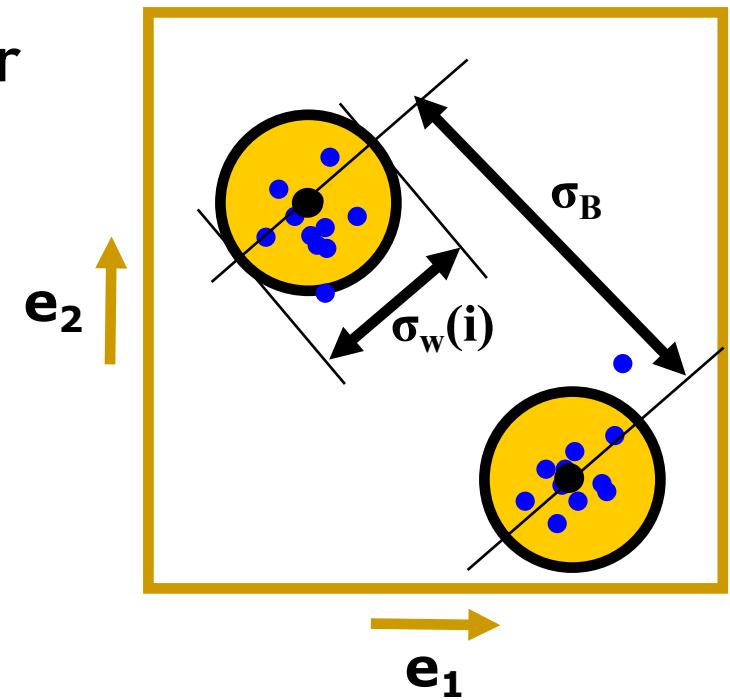
Clustering

Structure when:

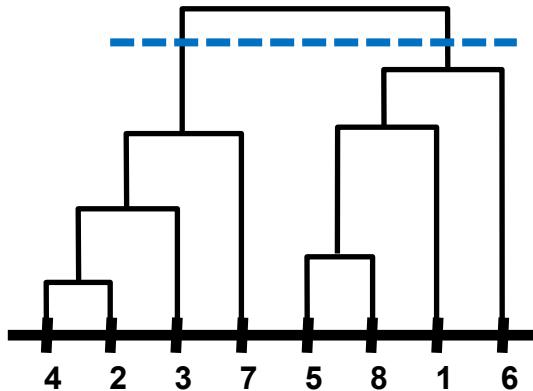
1. Samples within cluster resemble each other
(small within variance, $\sigma_w(i)$)
2. Clusters deviate from each other
(large between variance, σ_B)

Group samples such that:

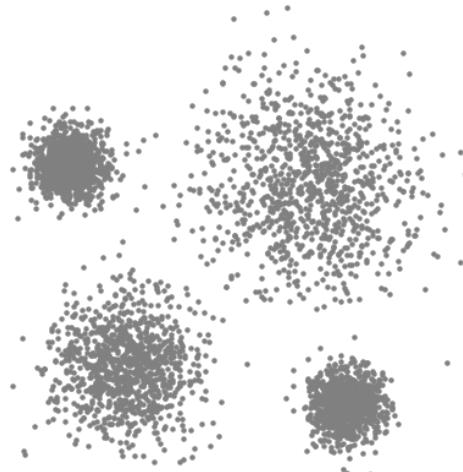
$$\min \left(\frac{\sum_{\text{clusters}} \sigma_w(i)}{\sigma_B} \right) \rightarrow \begin{array}{l} \sigma_w: \text{small} \\ \sigma_B: \text{large} \end{array}$$



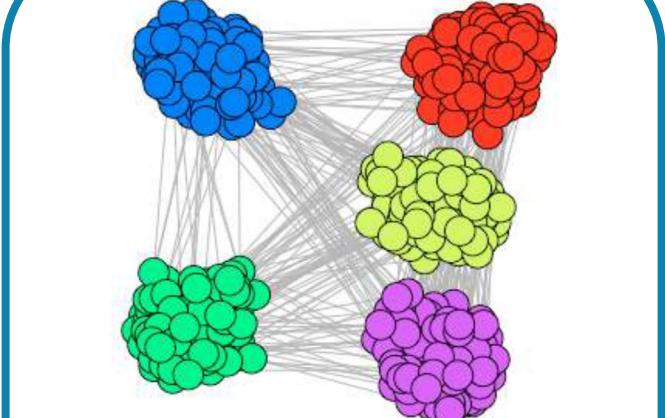
Many clustering approaches



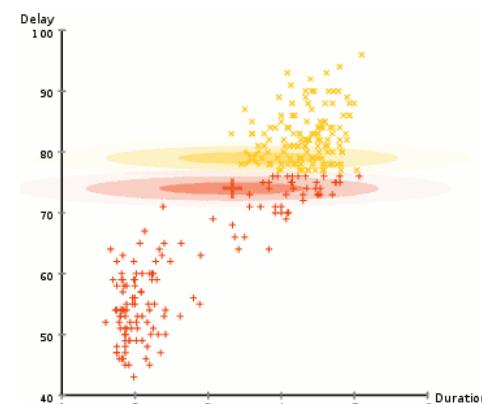
Hierarchical clustering



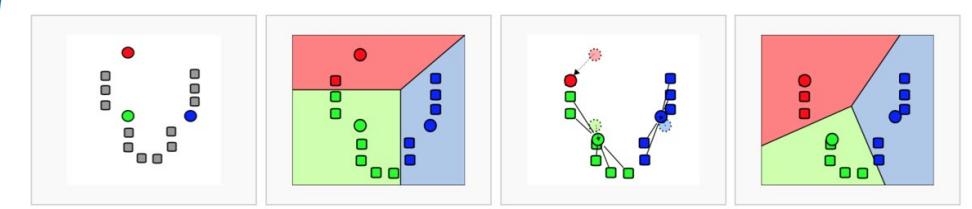
Mean shift clustering



Graph-based clustering



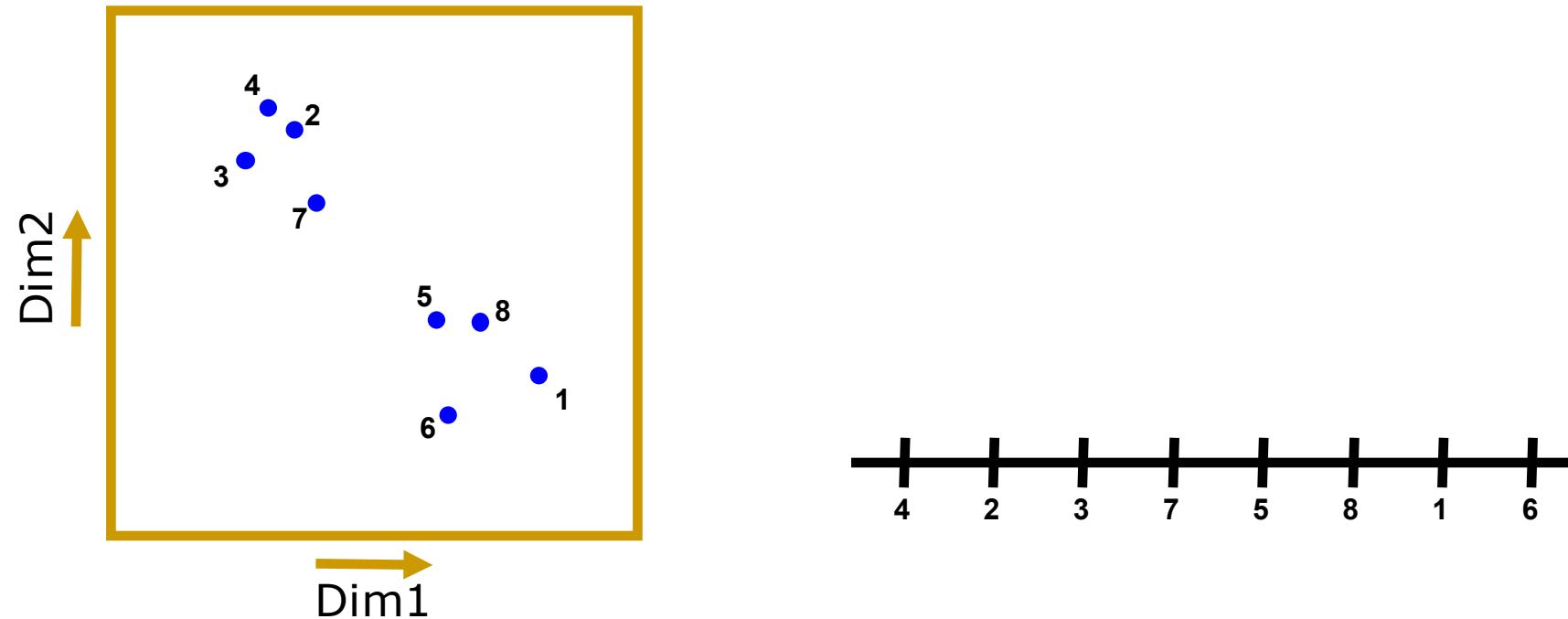
Gaussian mixture modeling



k -means clustering

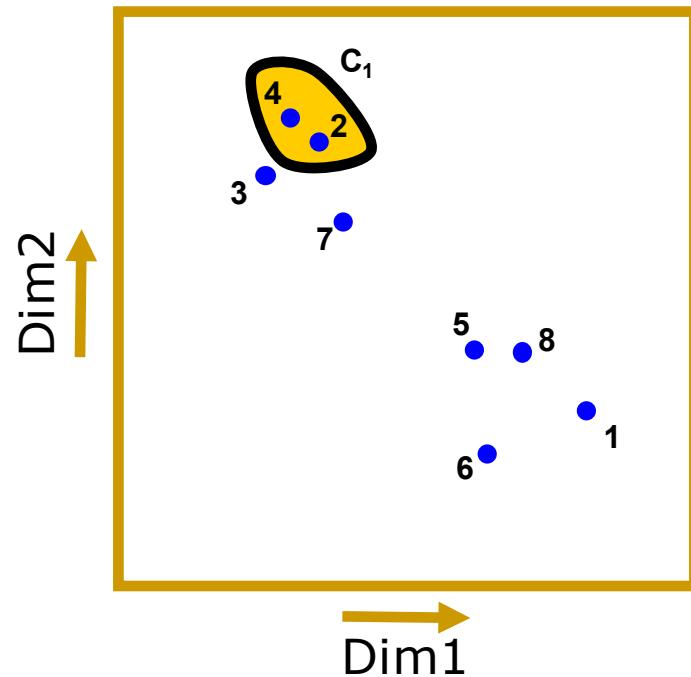
Herarchical clustering

Herarchical clustering

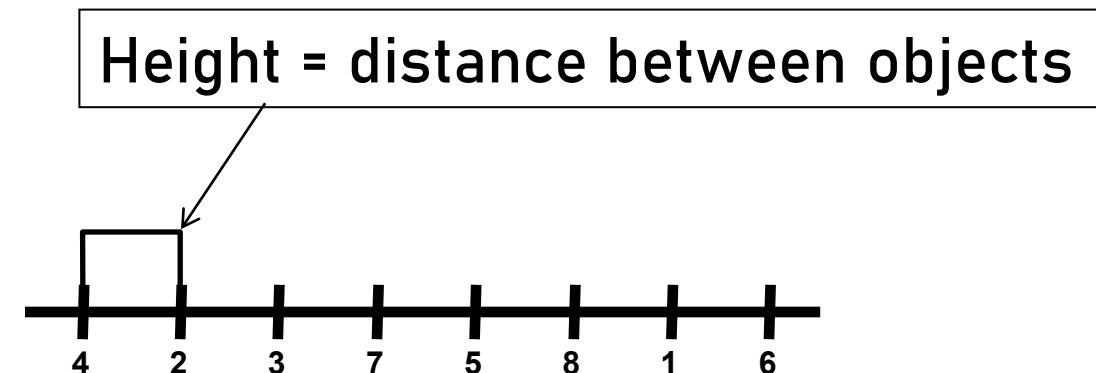


Find most similar objects (cells) and group them

Herarchical clustering



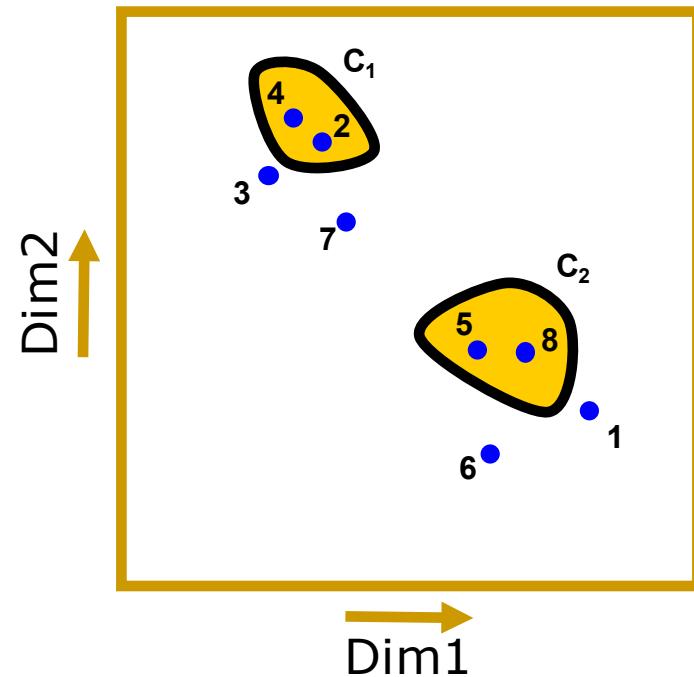
Dendrogram



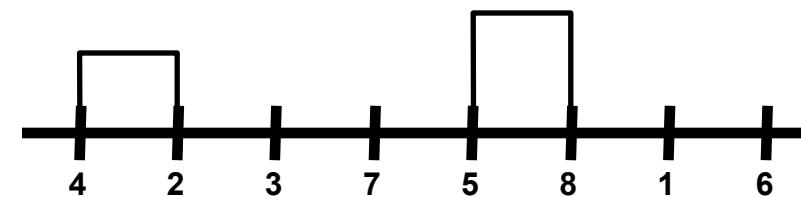
These are: objects 4 and 2

Again, find most similar objects (cells or clusters) and group them

Herarchical clustering



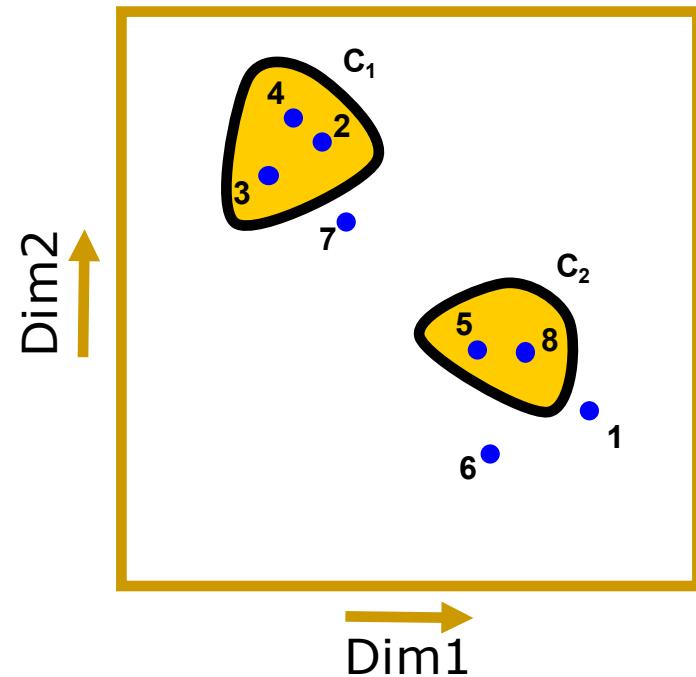
Dendrogram



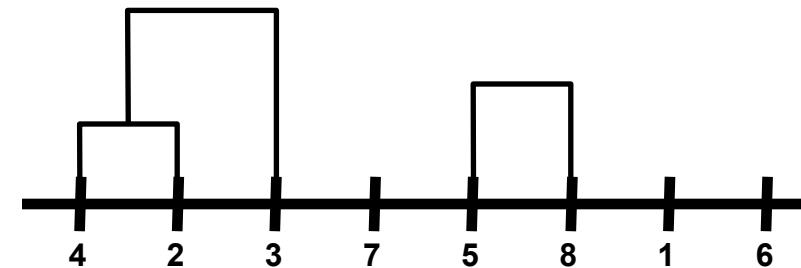
These are: objects 5 and 8

Repeat finding most similar objects (cells or clusters) and group them

Herarchical clustering

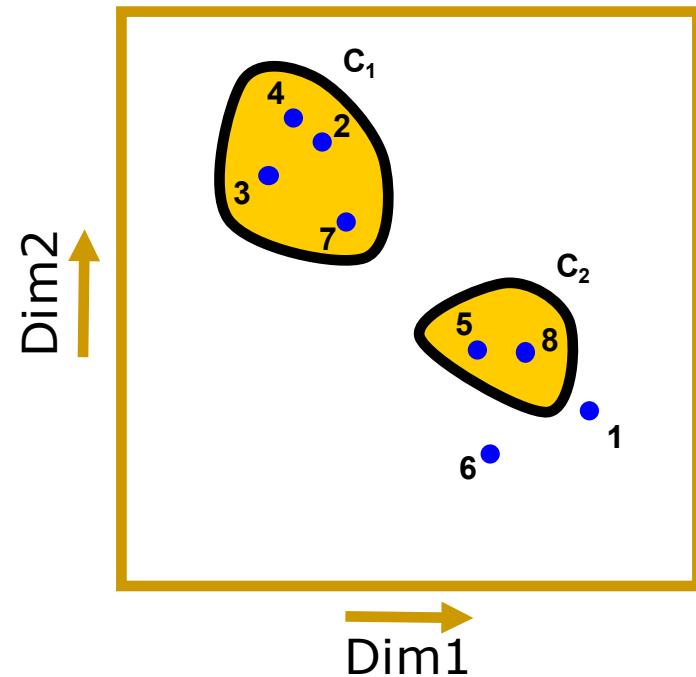


Dendrogram

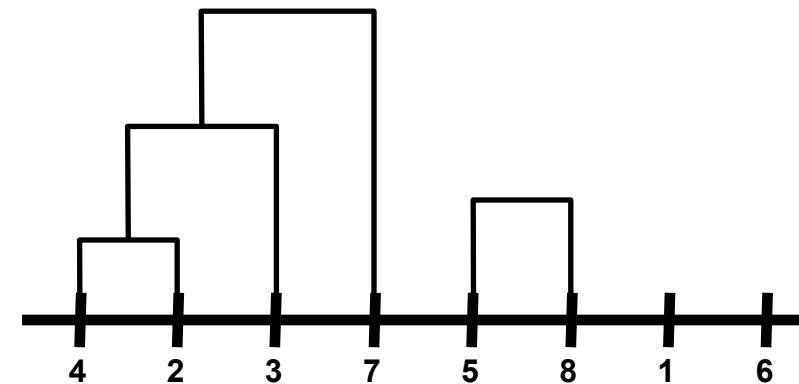


Join object 3 and cluster 1
Repeat process

Herarchical clustering

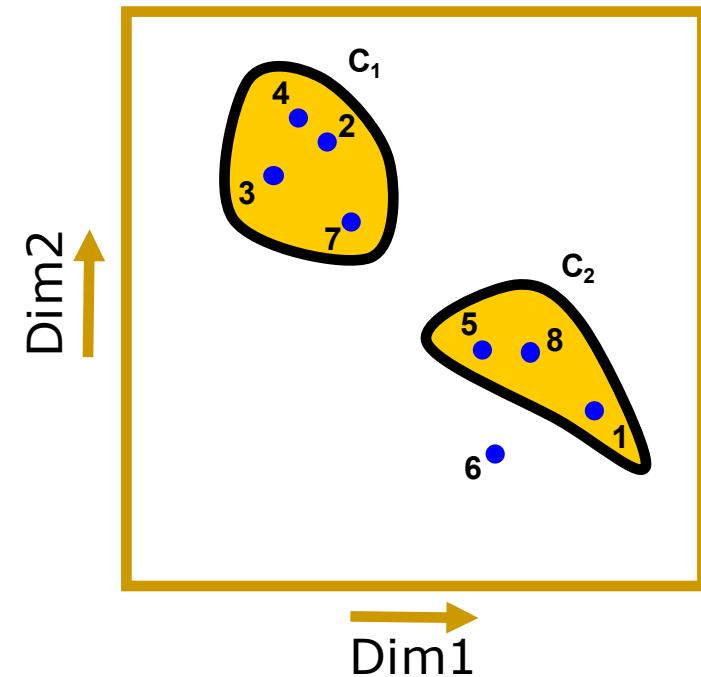


Dendrogram

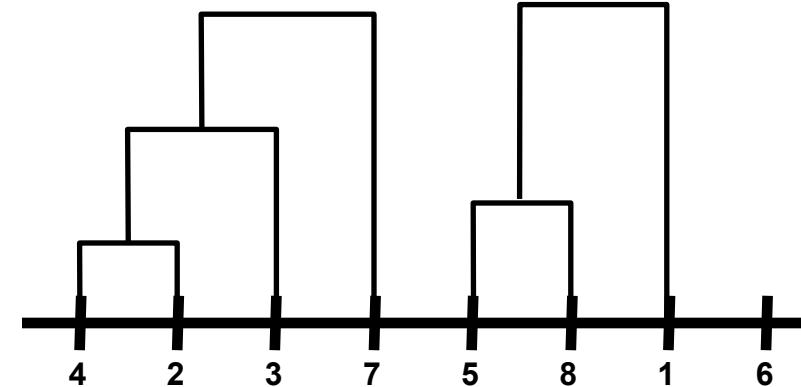


Join object 7 and cluster 1
Repeat process

Herarchical clustering

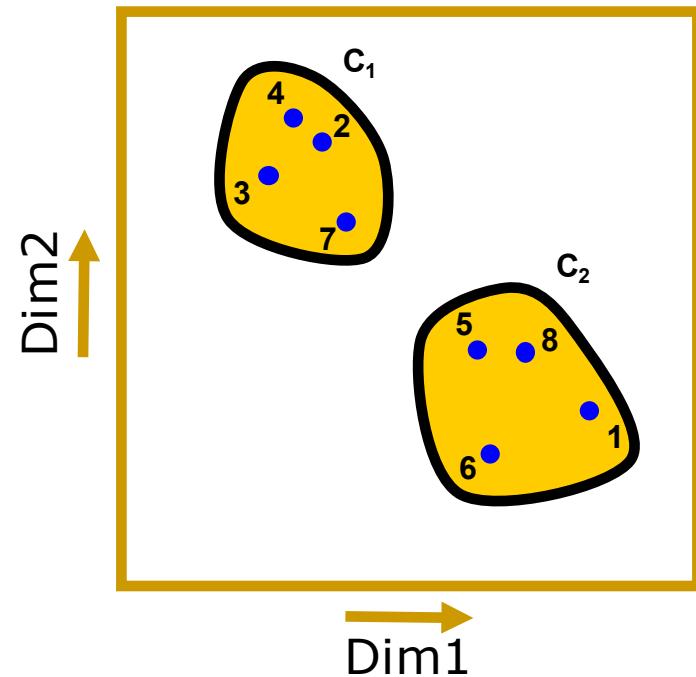


Dendrogram

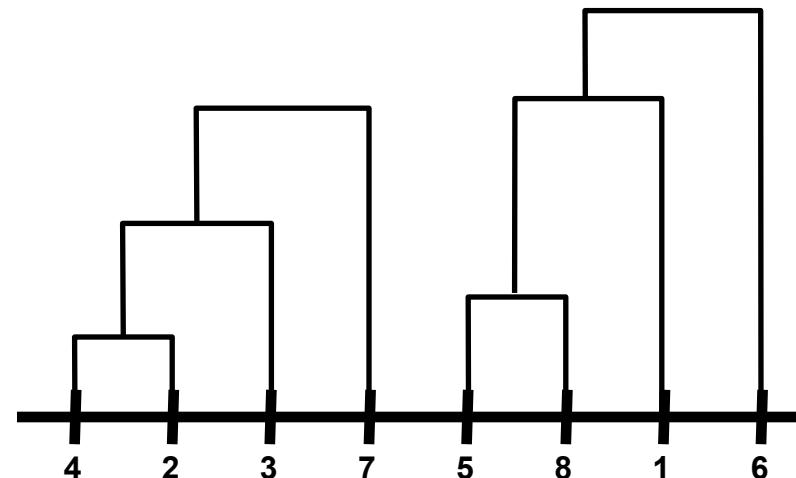


Join object 1 and cluster 2
Repeat process

Herarchical clustering

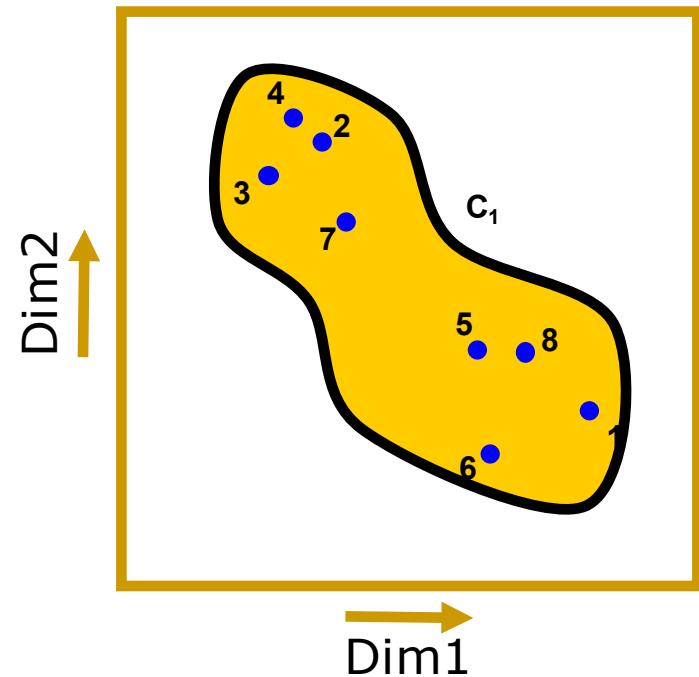


Dendrogram

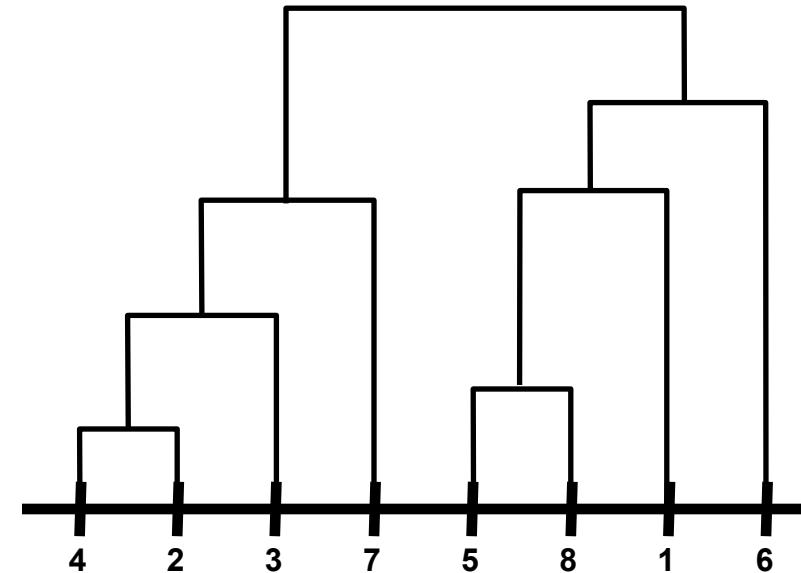


Join object 6 and cluster 2
Repeat process

Herarchical clustering

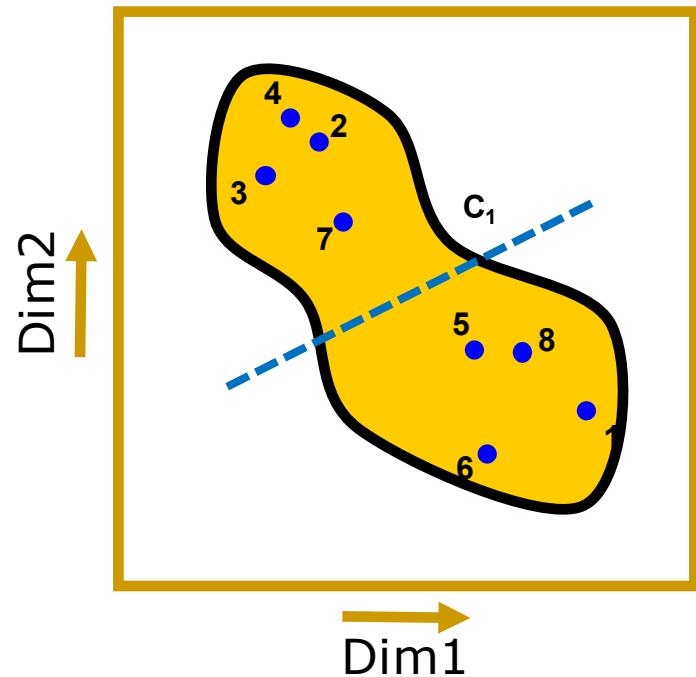


Dendrogram

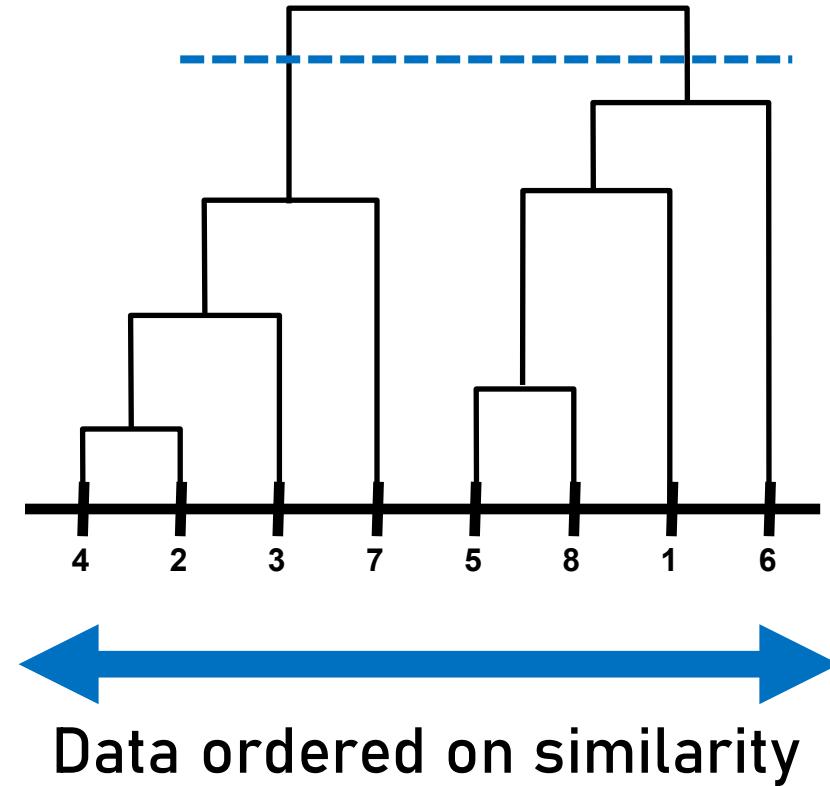


Join cluster 1 and cluster 2
All in one cluster: FINISHED!

Herarchical clustering



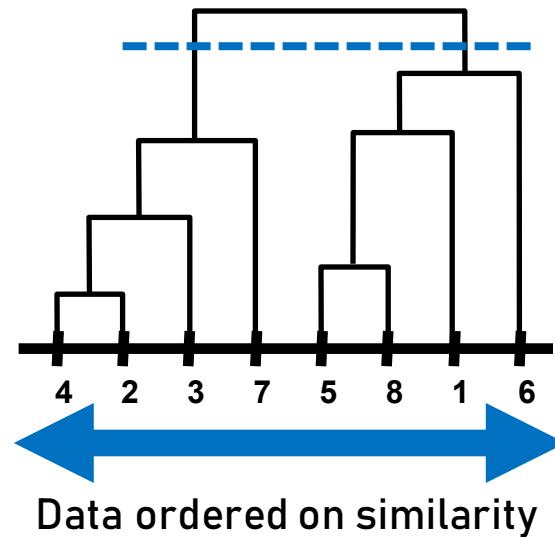
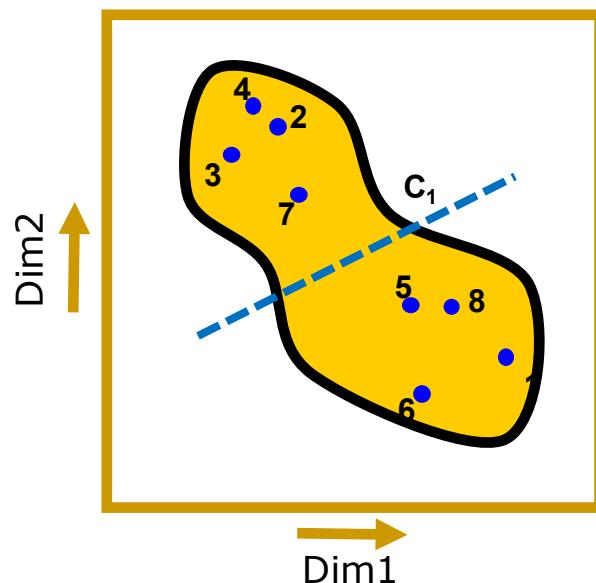
Dendrogram



Herarchical clustering

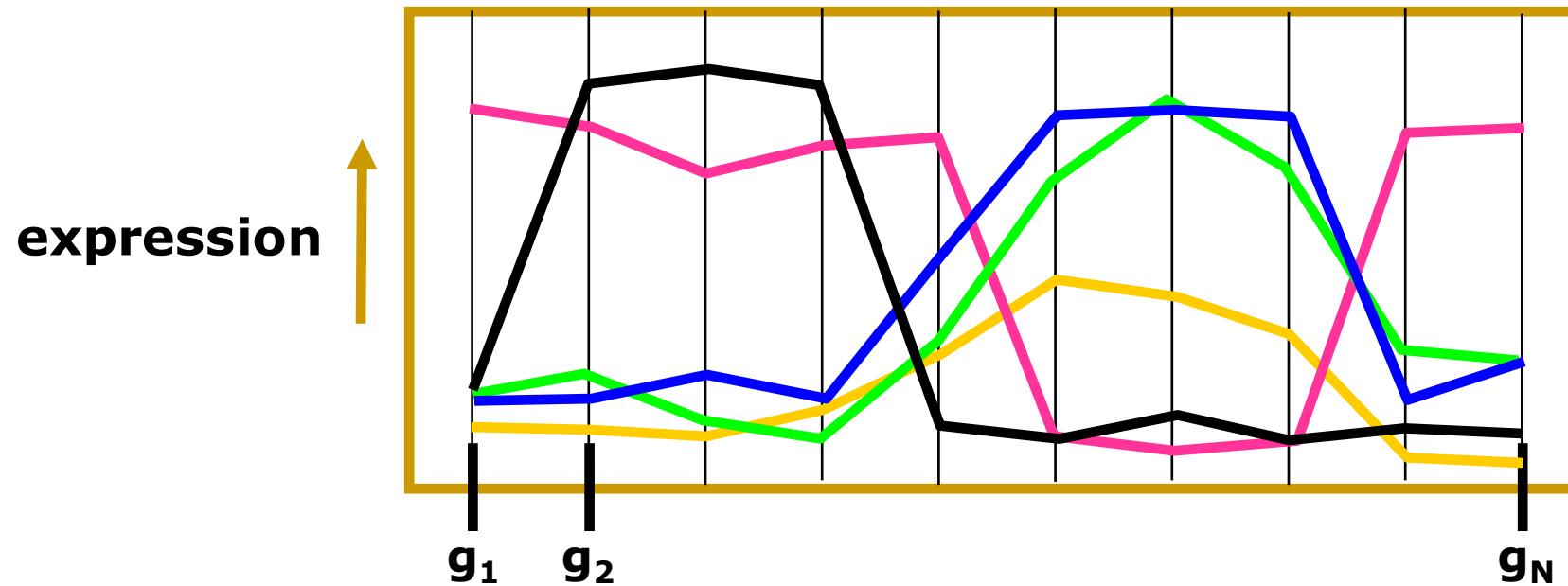
Need to know:

- Similarity between objects
- Similarity between clusters



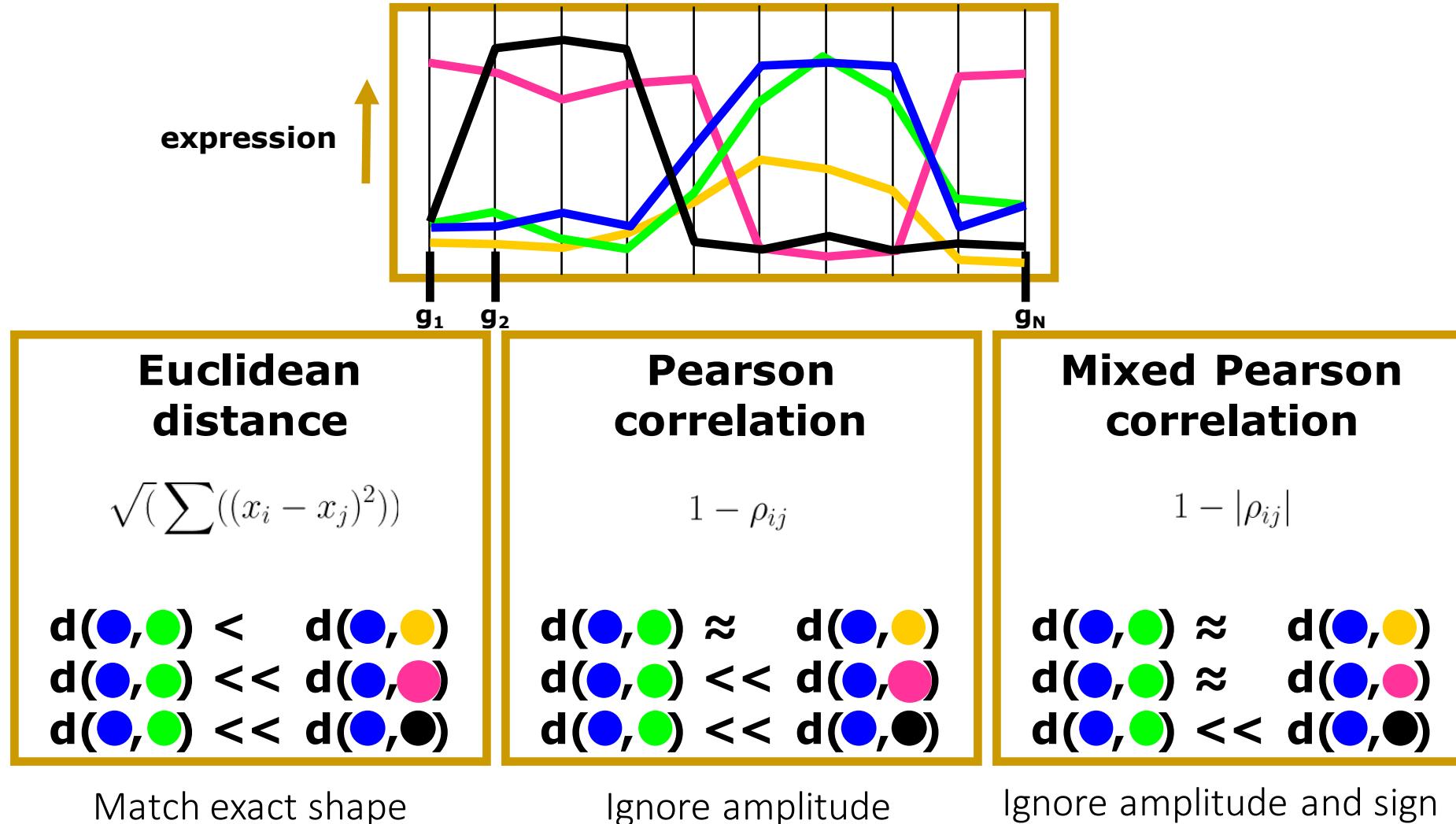
Herarchical clustering

Similarity between objects



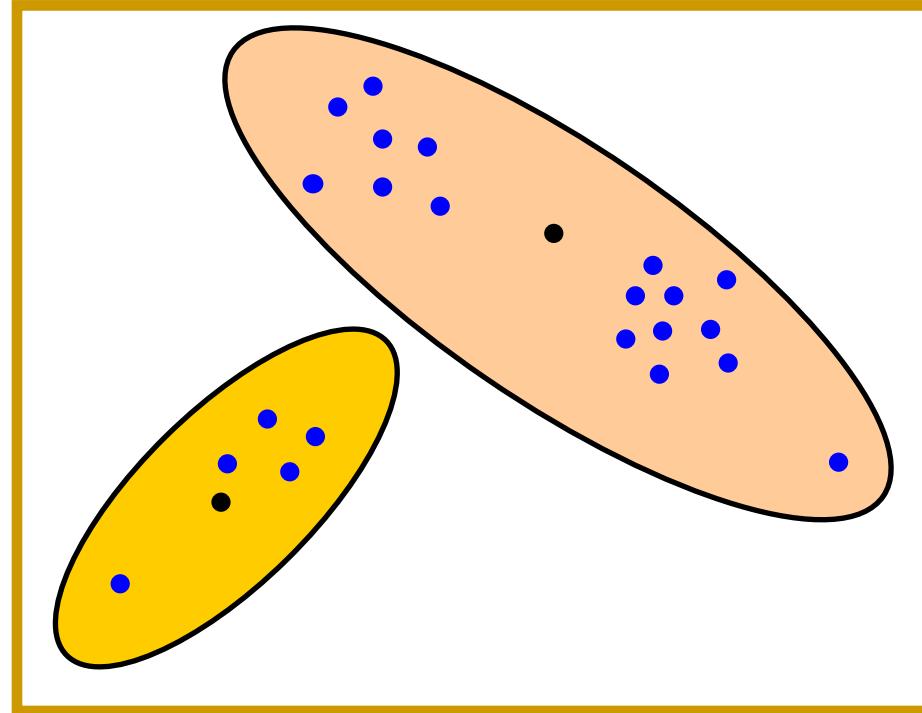
Herarchical clustering

Similarity between objects



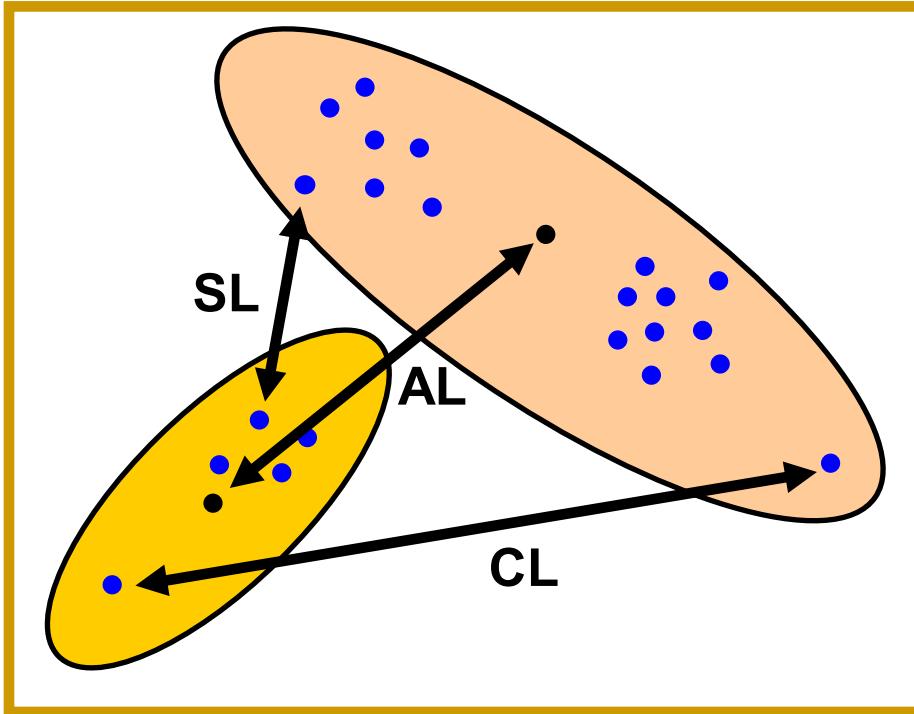
Hierarchical clustering

Similarity between clusters



Herarchical clustering

Similarity between clusters



Single linkage

Complete linkage

Average linkage

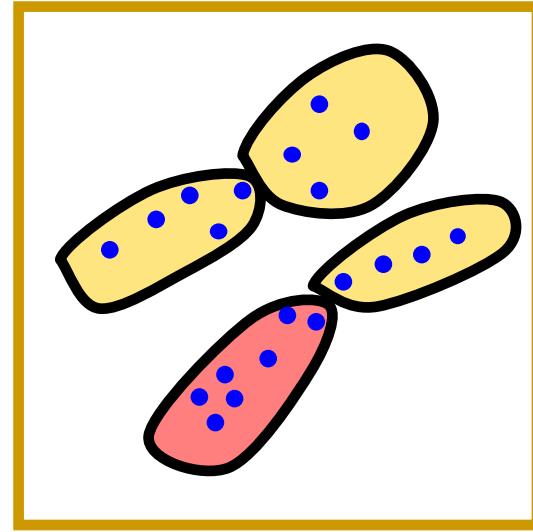
Closest objects

Furthest objects

Average similarity

Herarchical clustering

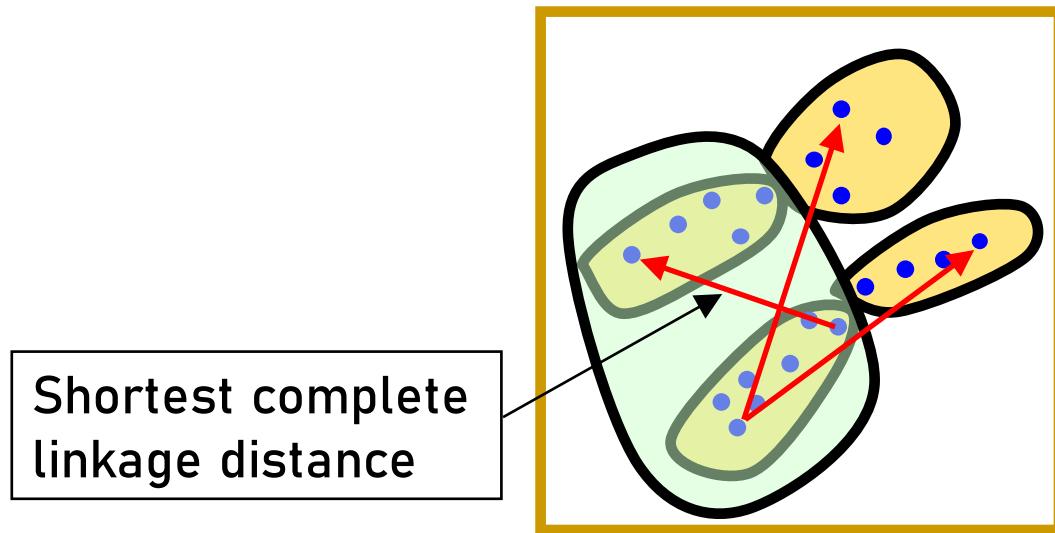
Similarity between clusters



Which cluster to merge with the red cluster when using **complete linkage**?

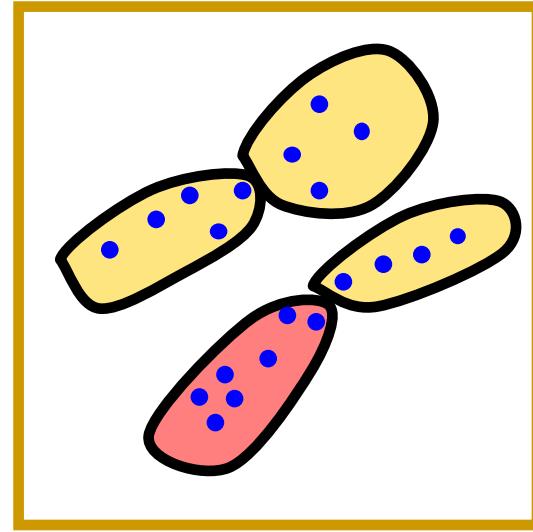
Herarchical clustering

Similarity between clusters



Herarchical clustering

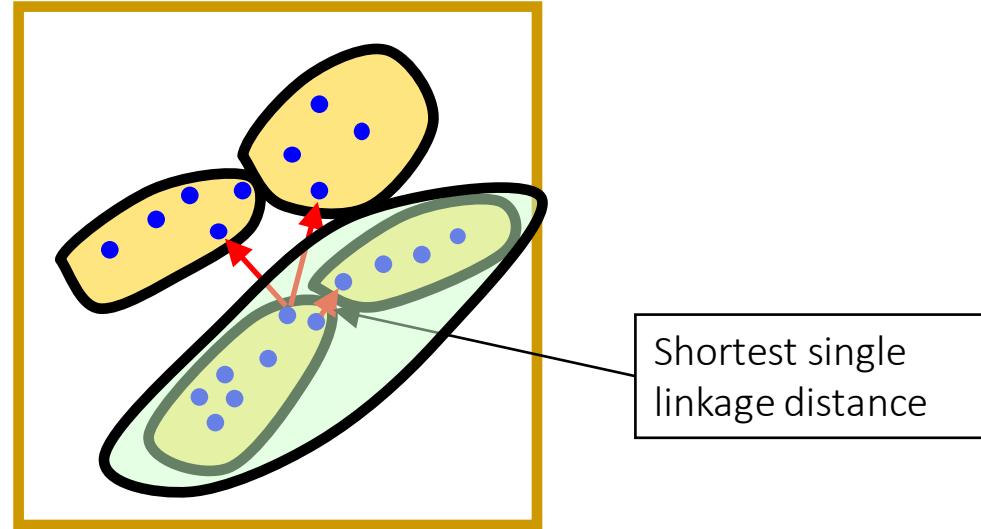
Similarity between clusters



Which cluster to merge with the red cluster when using **single linkage**?

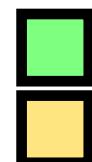
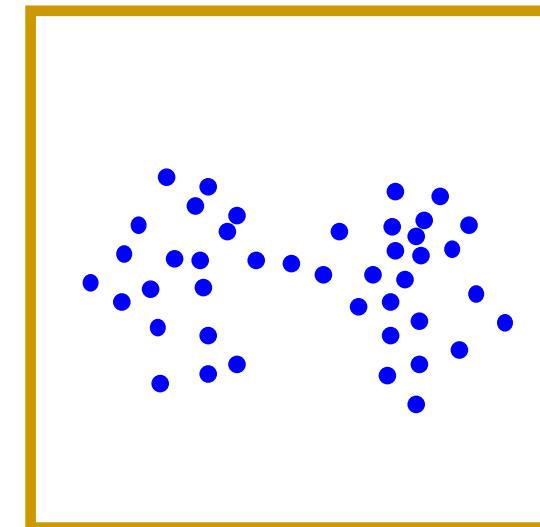
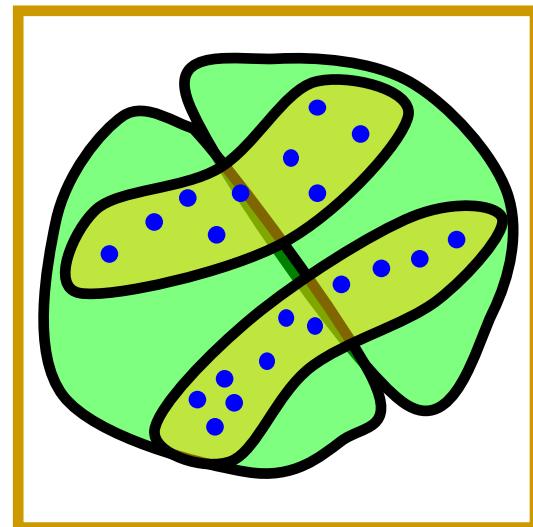
Herarchical clustering

Similarity between clusters



Herarchical clustering

Similarity between clusters

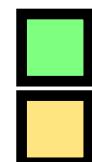
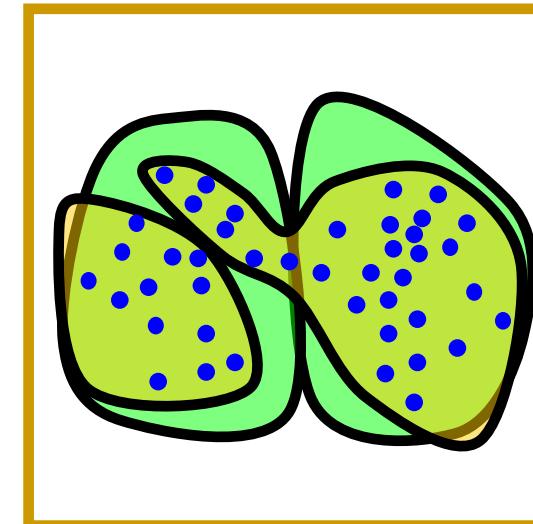
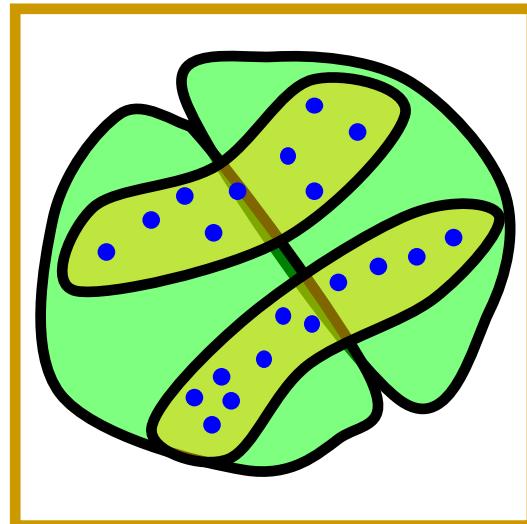


complete linkage
single linkage

Herarchical clustering

Similarity between clusters

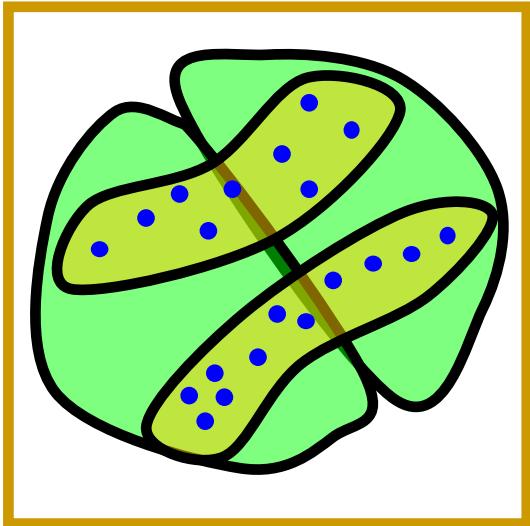
- Single linkage → long and “loose” clusters
- Complete linkage → compact clusters



complete linkage
single linkage

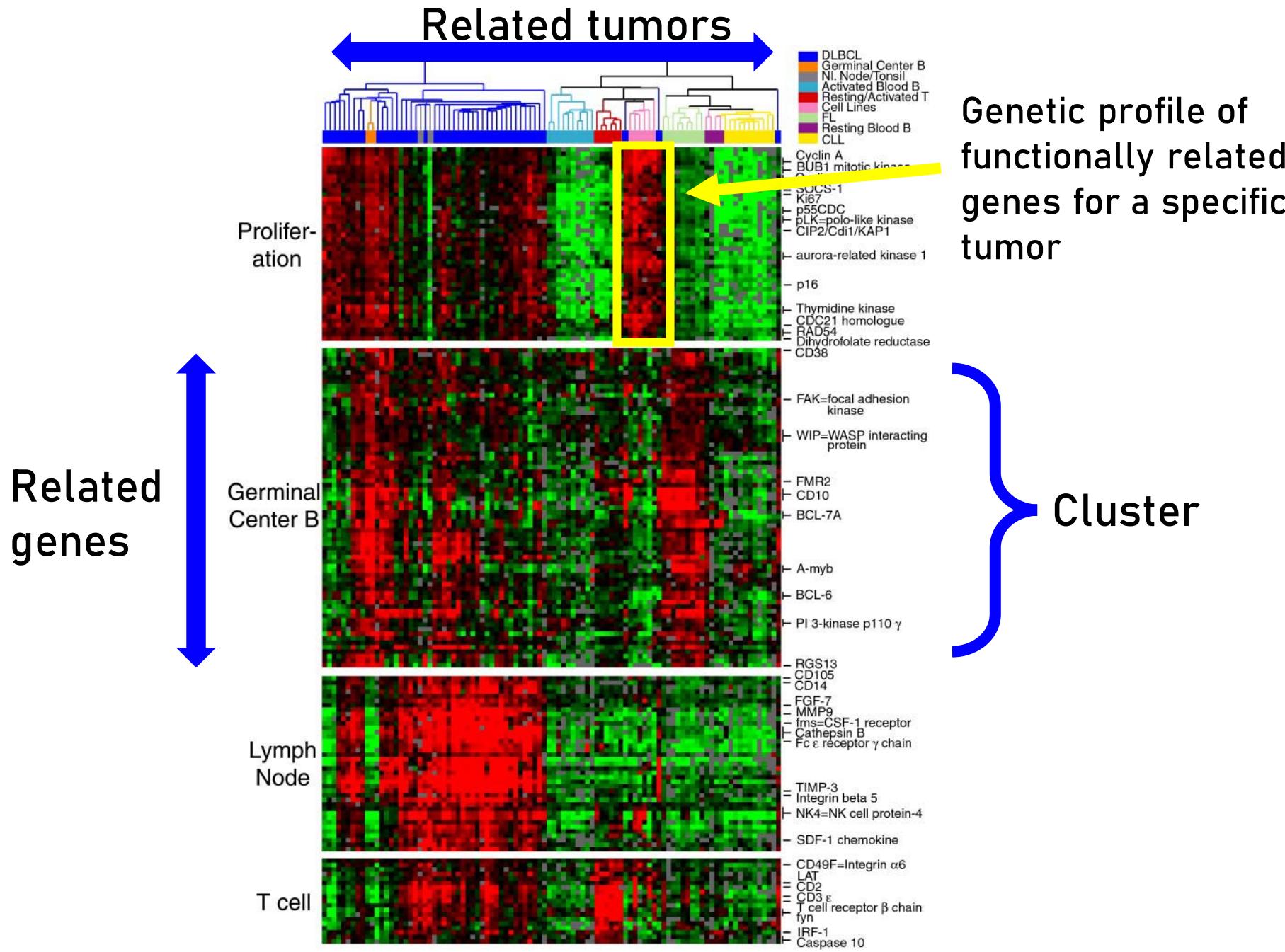
Herarchical clustering

Overview



complete linkage
single linkage

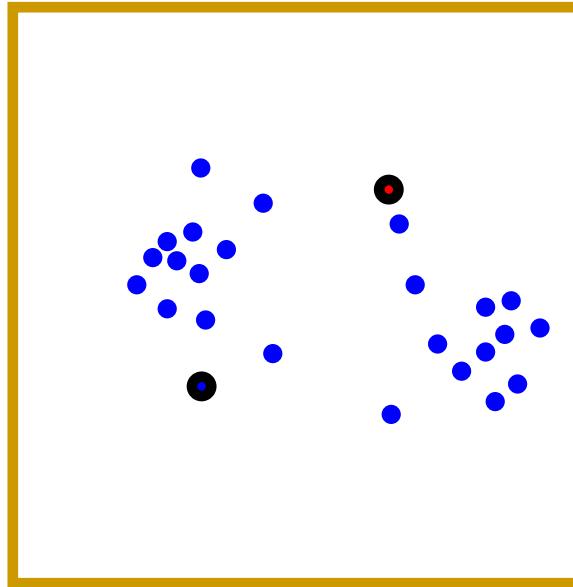
- Hierarchical clustering
 - Choice of distance measure
 - Choice of linkage type
- Distance measure
 - Euclidean
 - Correlation
- Linkage
 - Single
 - Average
 - Complete
- Number of clusters
 - Predefined or based on a cut-off in the dendrogram
 - Validate clustering!



Genetic profile of functionally related genes for a specific tumor

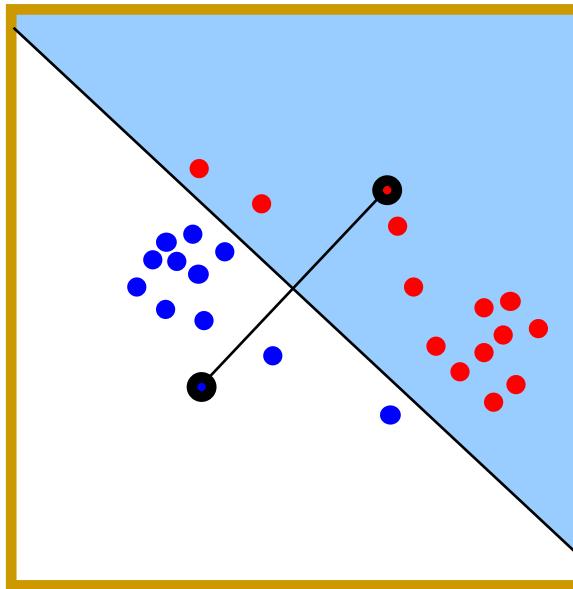
k -Means clustering

k -Means clustering



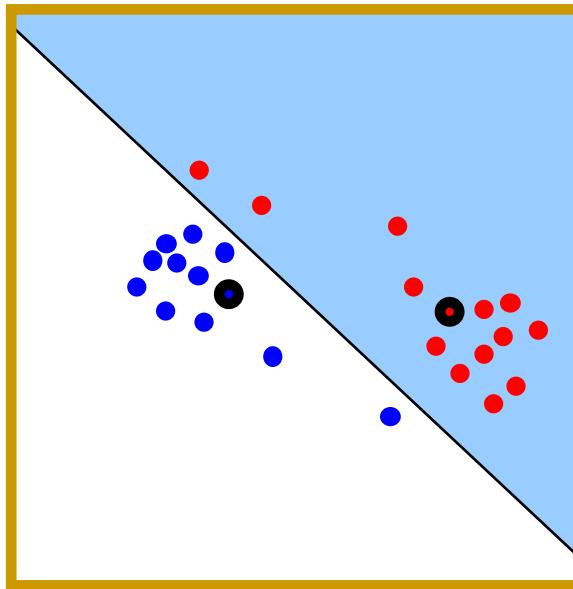
Choose randomly k prototypes

k -Means clustering



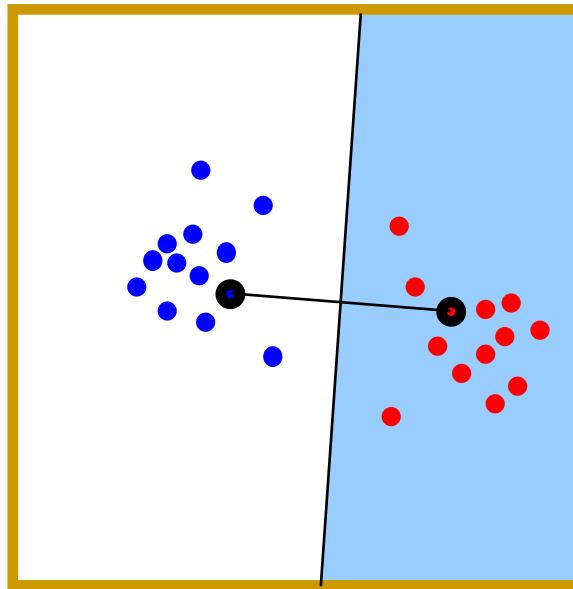
Assign objects to the closest prototype
Blue area: cluster 1
White area: cluster 2

k -Means clustering



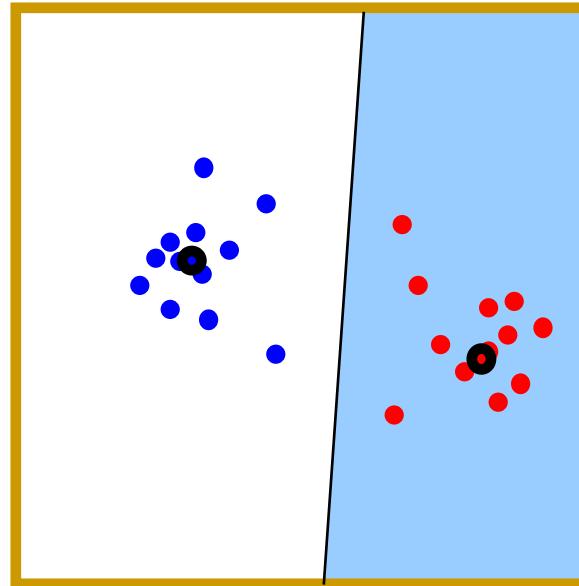
Calculate new cluster prototypes
By averaging objects

k -Means clustering



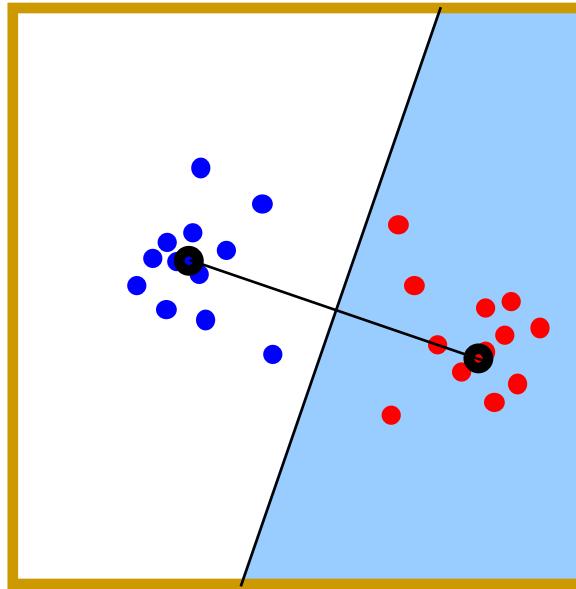
Re-assign objects to the closest prototype
Blue area: cluster 1
White area: cluster 2

k -Means clustering



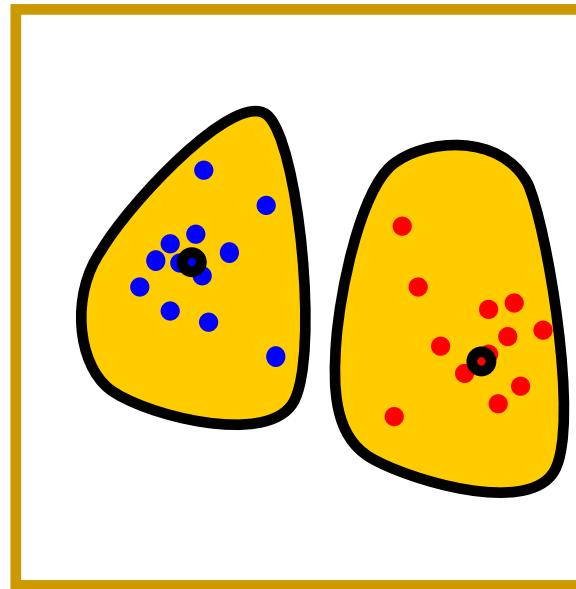
Re-calculate new cluster prototypes

k -Means clustering



Re-assign objects to the closest prototype
If no objects change cluster, then finished

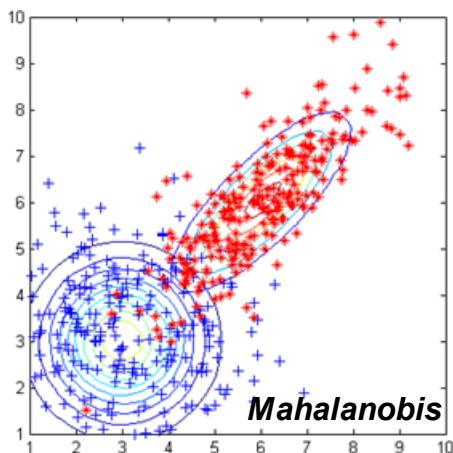
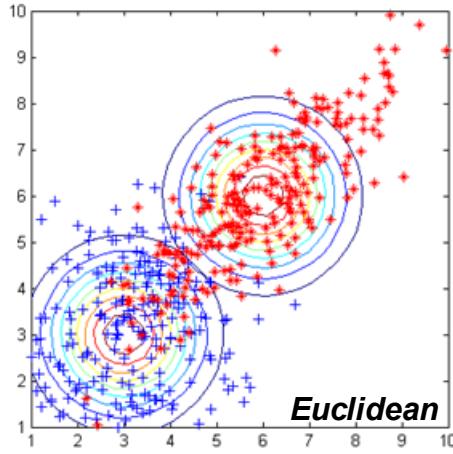
k -Means clustering



Establish clusters

k -Means clustering

Overview



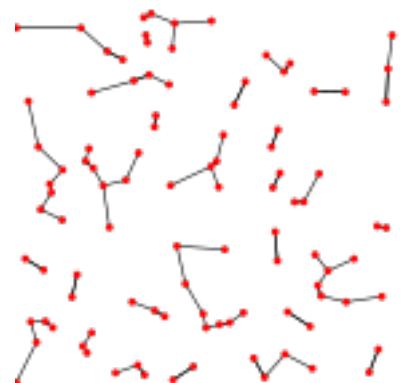
- **k -Means**
 - Fixed number of clusters (need to know a priori)
 - Choice of distance measure
 - Prototype choice
- **Distance measure**
 - Euclidean:
 - Mahalanobis:

Round clusters
Elongated clusters
- **Prototype choice**
 - Point
 - Line etc.
- **Number of clusters**
 - Predefined by k
 - Validate clustering!

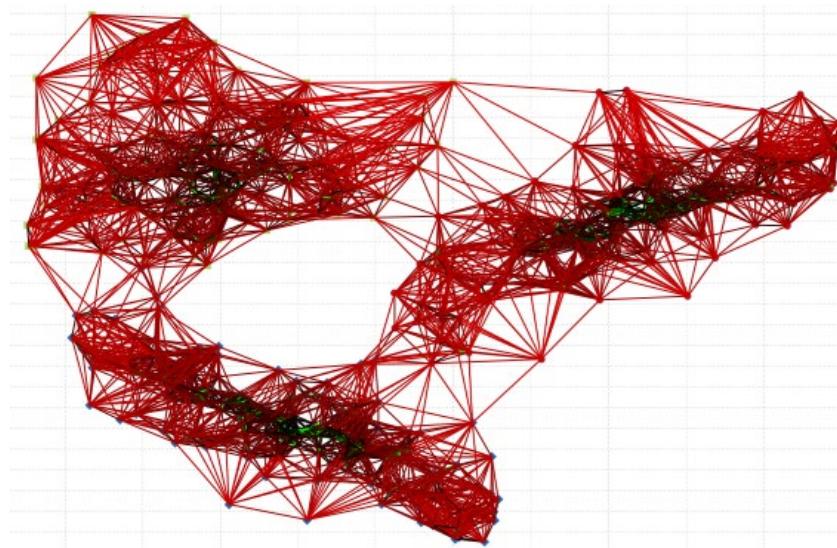
Graph-based clustering

Graph-based clustering

- k-NN graph: connect every node to its k-nearest neighbors
- Find densely connected components (communities)



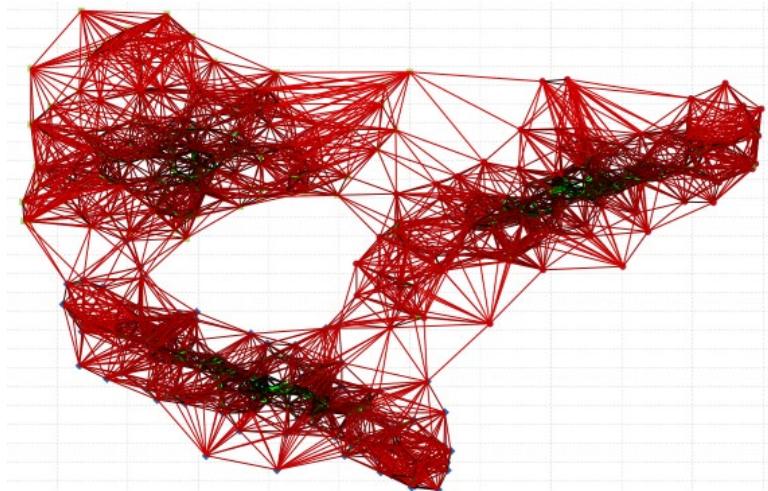
$k=1$



$k=20$

Graph-based clustering

- Maximize modularity score
 - Dense connections between nodes within clusters
 - Sparse connections between nodes in different clusters



Observed edges
in cluster c

Expected edges
in cluster c

$$H = \frac{1}{2m} \sum_c (e_c - \gamma \frac{K_c^2}{2m})$$

m = edges in the graph

e_c = edges in cluster c

$K_c = \sum_{n \in c} \text{degree}(n)$

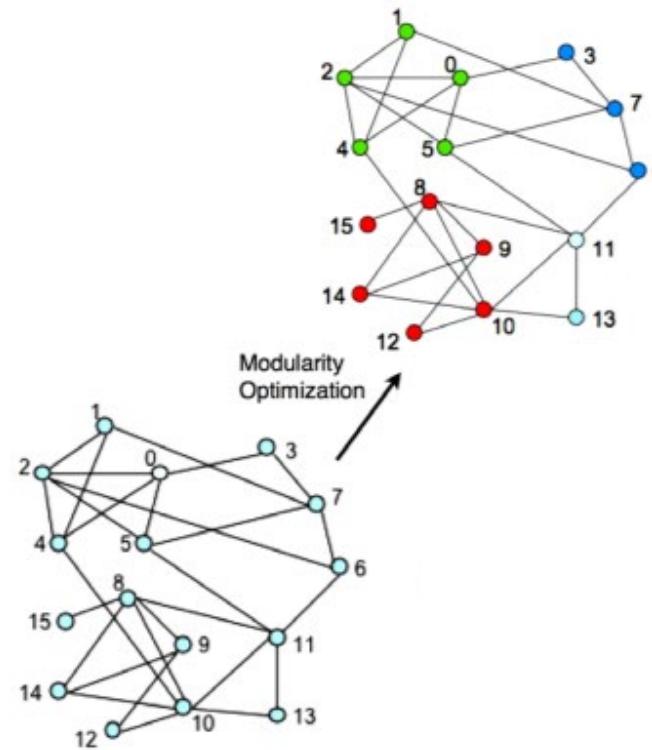
γ = resolution

Graph-based clustering

Louvain algorithm

Two steps

1. Local moving of nodes:
move node i to community
of neighbor j , if this
increases H
2. Aggregate nodes



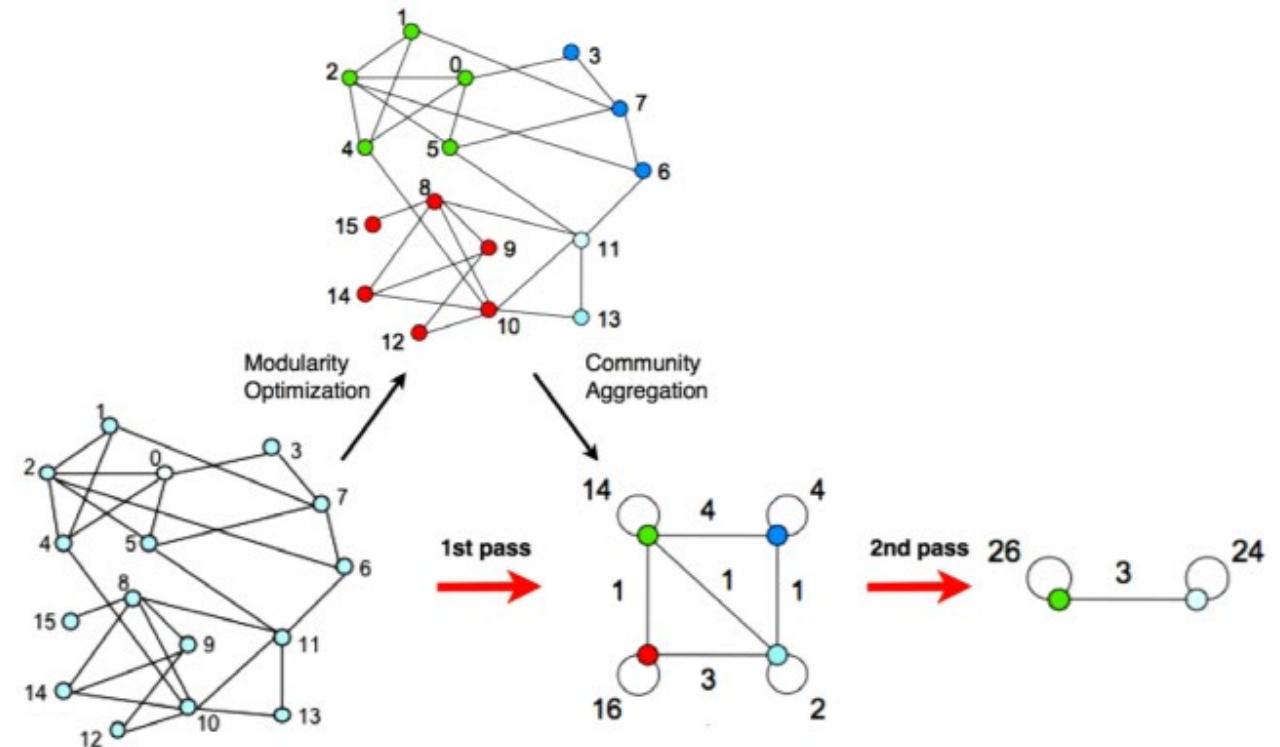
Graph-based clustering

Louvain algorithm

Two steps

1. Local moving of nodes:
move node i to community
of neighbor j , if this
increases H
2. Aggregate nodes

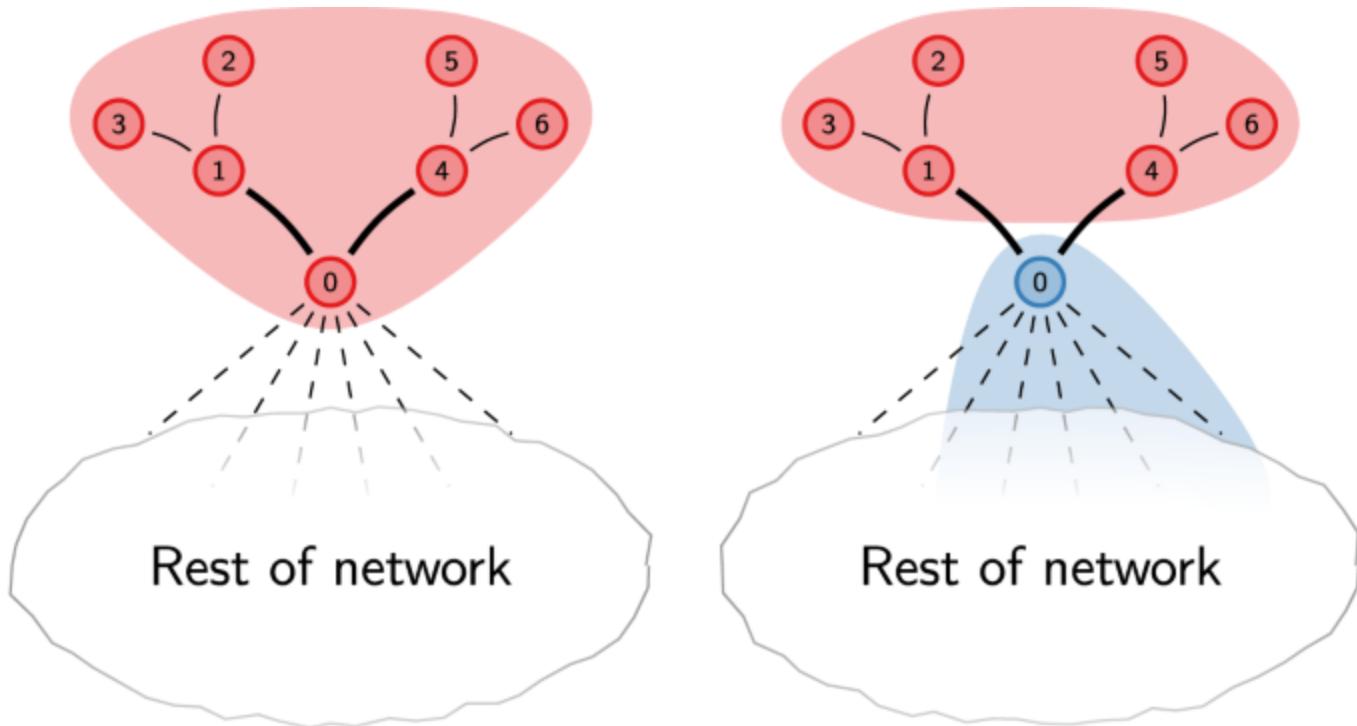
Iterate until no more changes



Graph-based clustering

Louvain algorithm

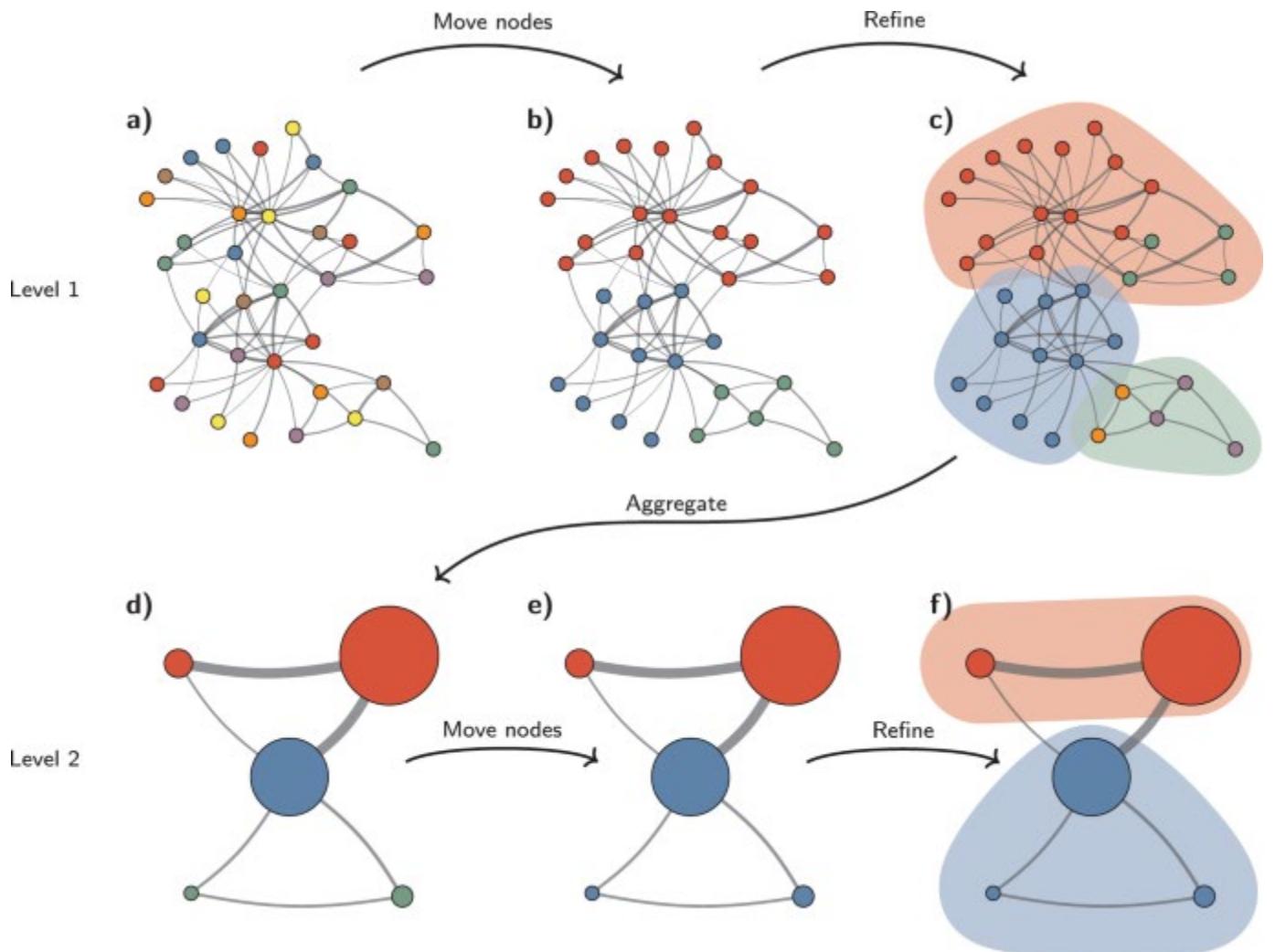
- During the ‘moving step’, nodes can become internally disconnected
- Nodes 1-6 still locally optimal assigned



Graph-based clustering

Leiden algorithm

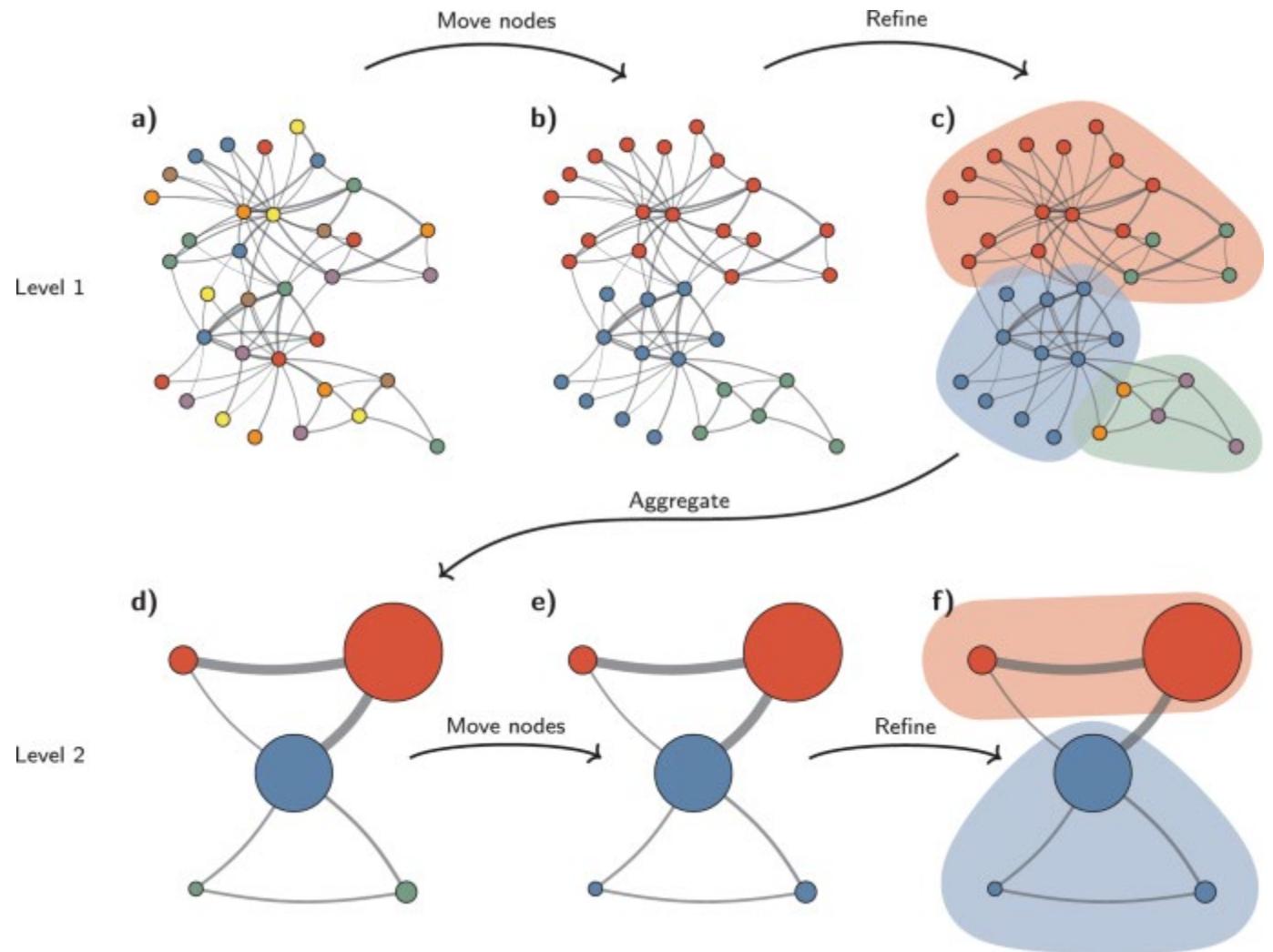
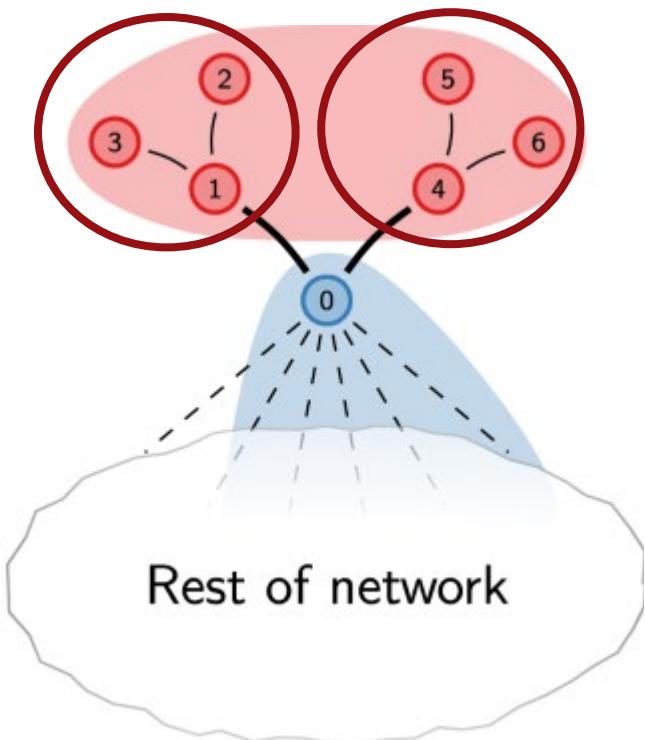
- Solution: add refinement step



Graph-based clustering

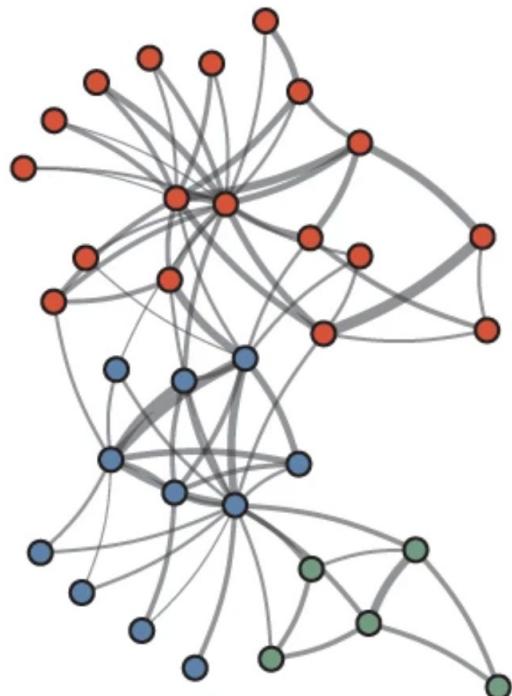
Leiden algorithm

- Solution: add refinement step



Graph-based clustering

Overview

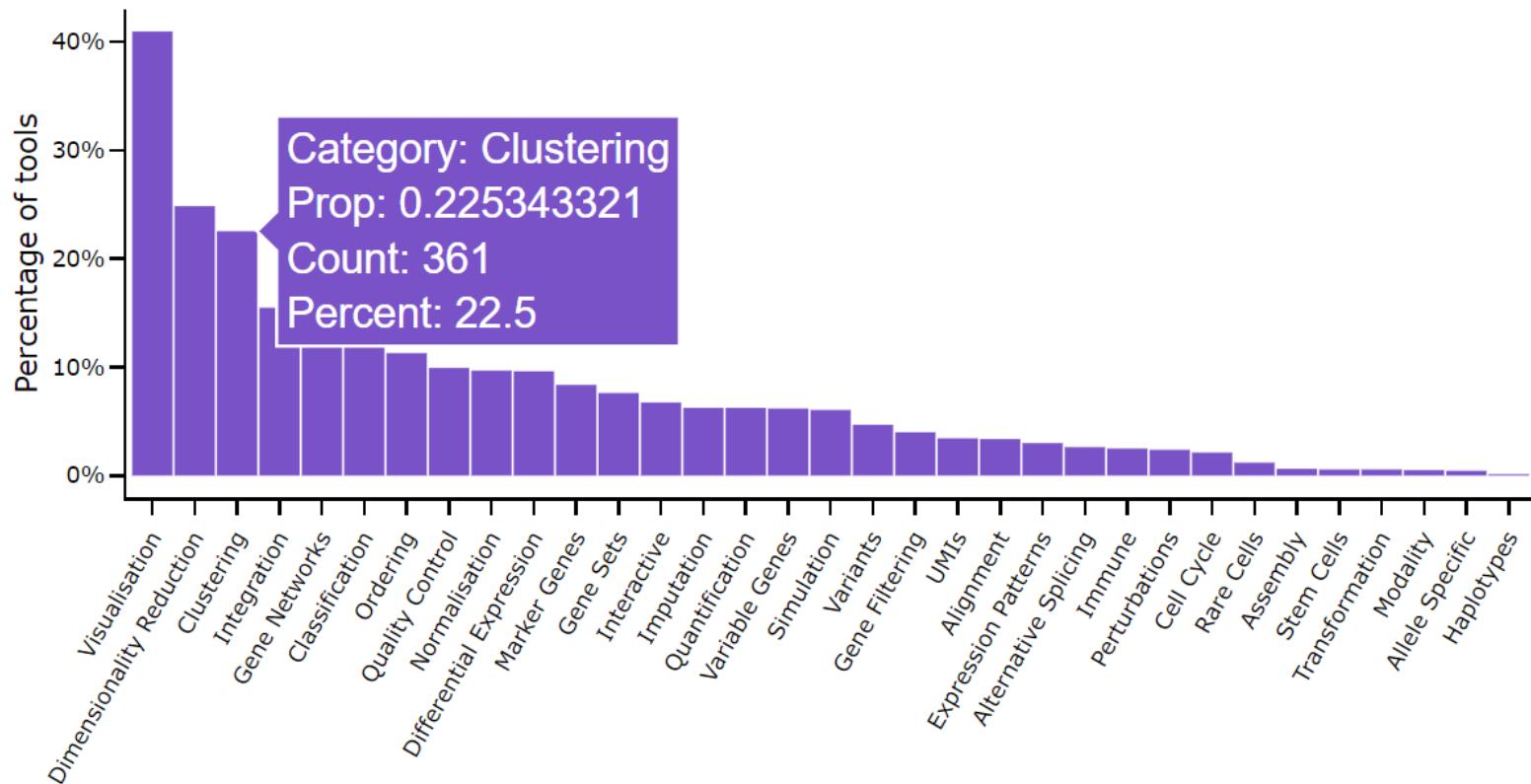


- **Graph-based clustering**
 - Number of neighbors when constructing the graph
 - Resolution parameters
- **Resolution**
 - High → less clusters
 - Low → more clusters
- **Number of clusters**
 - Determined using resolution parameter
 - Validate clustering!

Outline

- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- scRNA-seq clustering
- Annotating clusters
- Cluster validation
- Challenges & outlook

>300 scRNA-seq clustering methods available



scRNA-seq clustering methods

Name	Year	Method type	Strengths	Limitations
scanpy ⁴	2018	PCA+graph-based	Very scalable	May not be accurate for small data sets
Seurat (latest) ³	2016			
PhenoGraph ³²	2015			
SC3 (REF. ²²)	2017	PCA+k-means	High accuracy through consensus, provides estimation of k	High complexity, not scalable
SIMLR ²⁴	2017	Data-driven dimensionality reduction+k-means	Concurrent training of the distance metric improves sensitivity in noisy data sets	Adjusting the distance metric to make cells fit the clusters may artificially inflate quality measures
CIDR ²⁵	2017	PCA+hierarchical	Implicitly imputes dropouts when calculating distances	
GiniClust ⁷⁵	2016	DBSCAN	Sensitive to rare cell types	Not effective for the detection of large clusters
pcaReduce ²⁷	2016	PCA+k-means+hierarchical	Provides hierarchy of solutions	Very stochastic, does not provide a stable result
Tasic et al. ²⁸	2016	PCA+hierarchical	Cross validation used to perform fuzzy clustering	High complexity, no software package available
TSCAN ⁴¹	2016	PCA+Gaussian mixture model	Combines clustering and pseudotime analysis	Assumes clusters follow multivariate normal distribution
mpath ⁴⁵	2016	Hierarchical	Combines clustering and pseudotime analysis	Uses empirically defined thresholds and a priori knowledge
BackSPIN ²⁶	2015	Biclustering (hierarchical)	Multiple rounds of feature selection improve clustering resolution	Tends to over-partition the data
RaceID ²³ , RaceID2 (REF. ¹¹⁵), RaceID3	2015	k-Means	Detects rare cell types, provides estimation of k	Performs poorly when there are no rare cell types
SINCERA ⁵	2015	Hierarchical	Method is intuitively easy to understand	Simple hierarchical clustering is used, may not be appropriate for very noisy data
SNN-Clip ⁸⁰	2015	Graph-based	Provides estimation of k	High complexity, not scalable

How to compare different cluster labels?

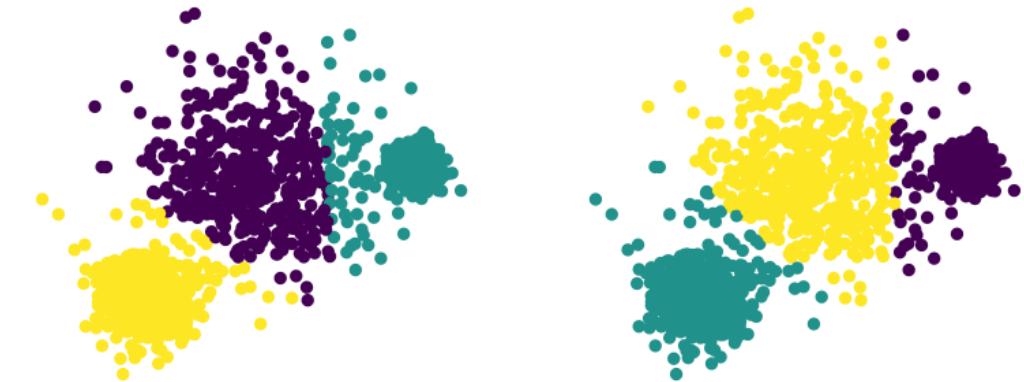
Adjusted Rand Index (ARI)

Measure of the similarity between two data clusterings

Given a set S of n elements, and two groupings or partitions of these elements

$$X = \{X_1, X_2, \dots, X_r\}, \quad Y = \{Y_1, Y_2, \dots, Y_s\}$$

$X \setminus Y$	Y_1	Y_2	\dots	Y_s	sums
X_1	n_{11}	n_{12}	\dots	n_{1s}	a_1
X_2	n_{21}	n_{22}	\dots	n_{2s}	a_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
X_r	n_{r1}	n_{r2}	\dots	n_{rs}	a_r
sums	b_1	b_2	\dots	b_s	



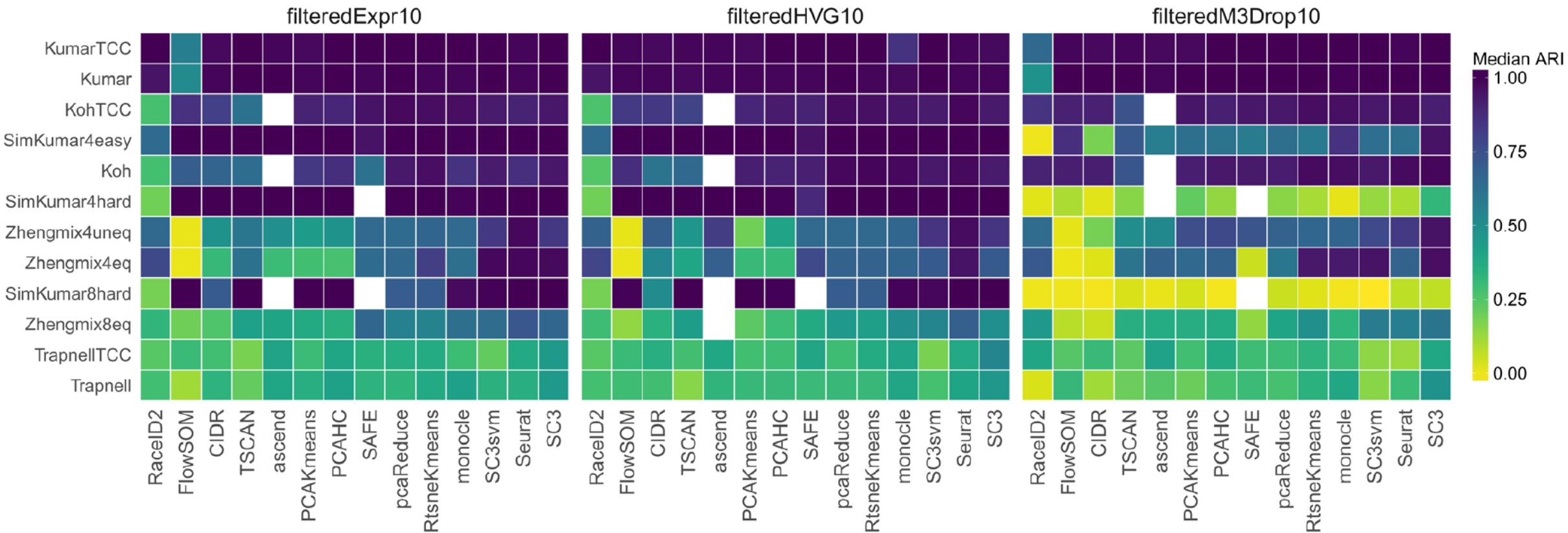
$$ARI = \frac{\underbrace{\sum_{ij} \binom{n_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}_{\text{Index}}}{\underbrace{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}_{\text{Max index}}}$$

Expected index

Expected index

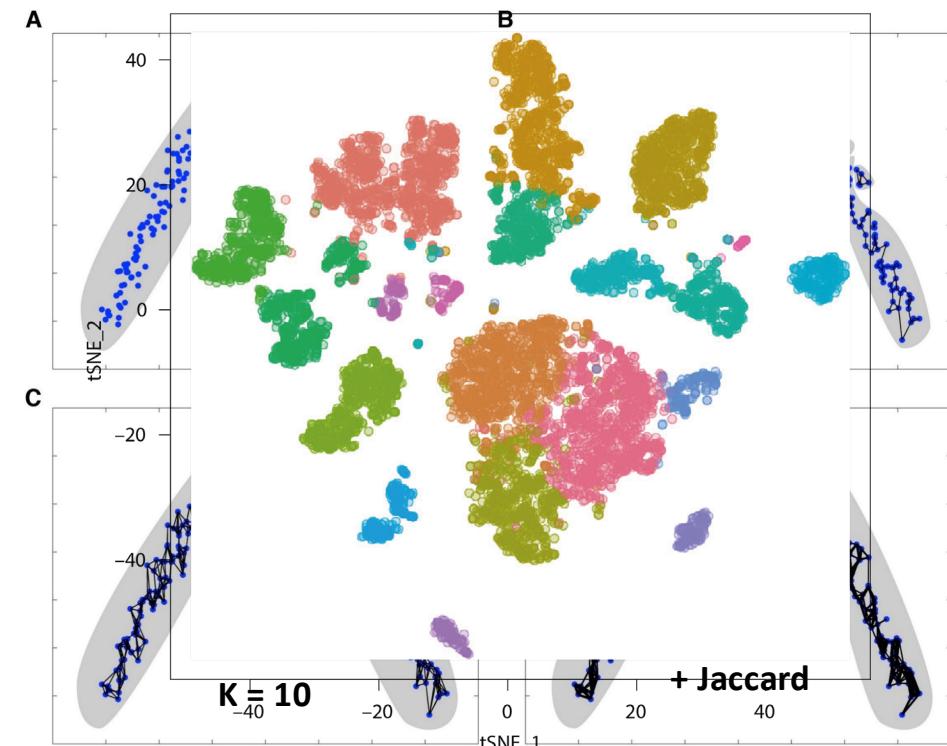
Max index

Benchmarking scRNA-seq clustering methods



Standard clustering approach

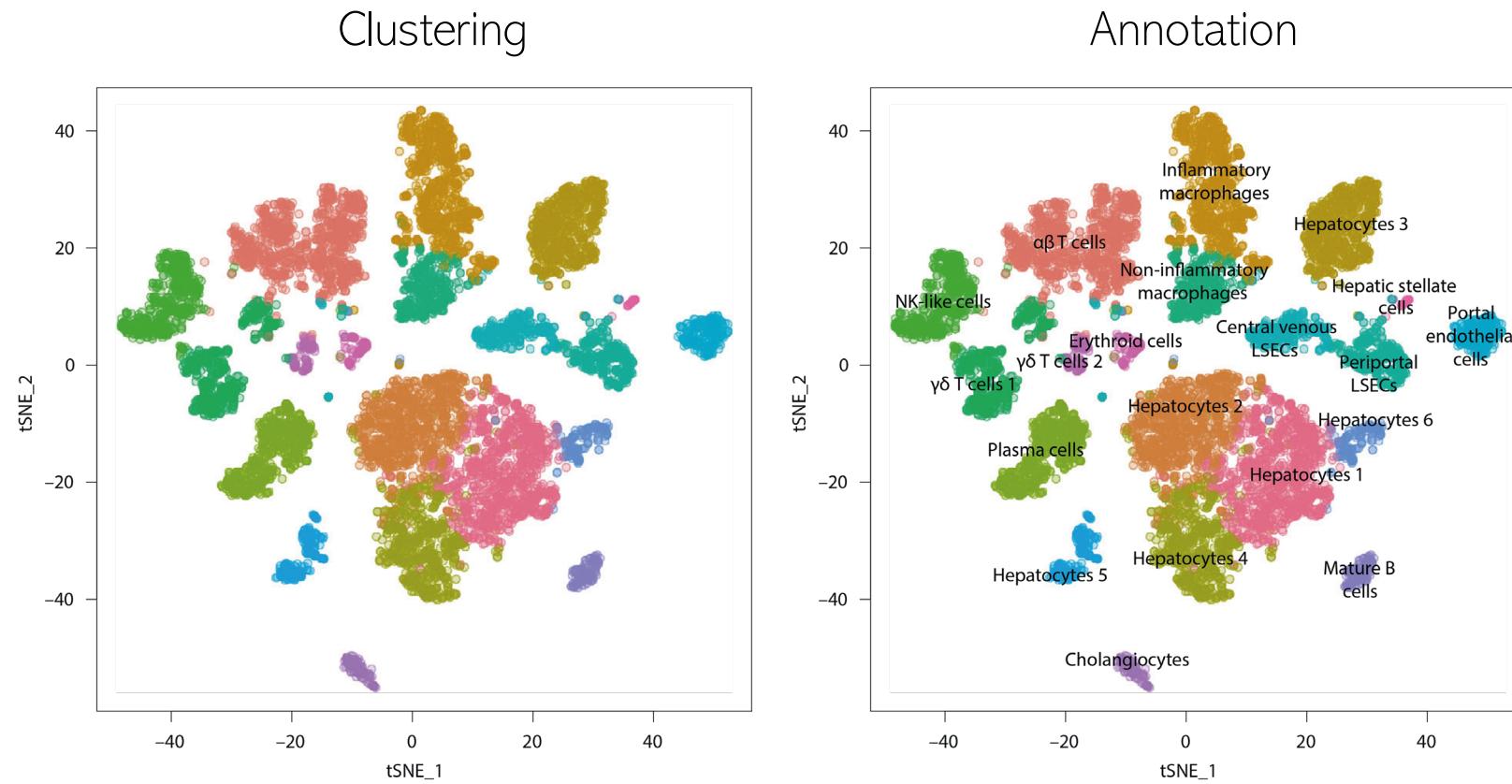
1. Select highly variable genes (~1000-5000 genes)
2. Reduce dimensions using PCA (~30-50 dimensions)
3. Construct kNN graph (~15-20 neighbors)
4. Louvain/Leiden community detection



Outline

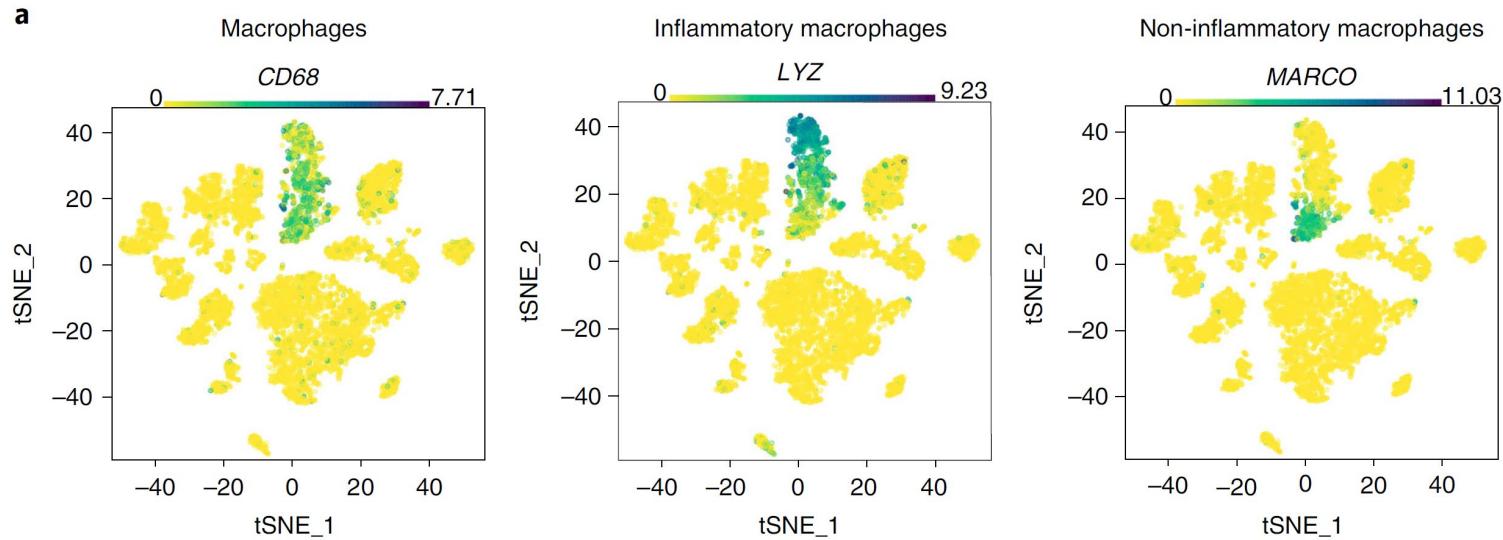
- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- scRNA-seq clustering
- Annotating clusters
- Cluster validation
- Challenges & outlook

From clusters to annotations

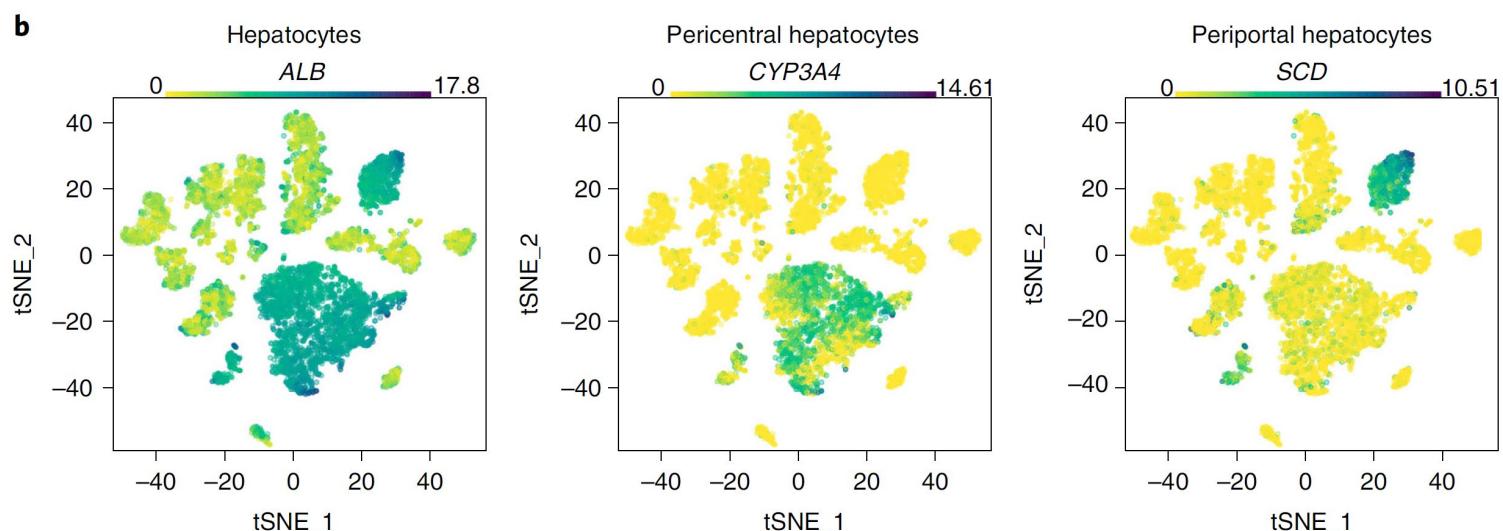


Gene expression overlay

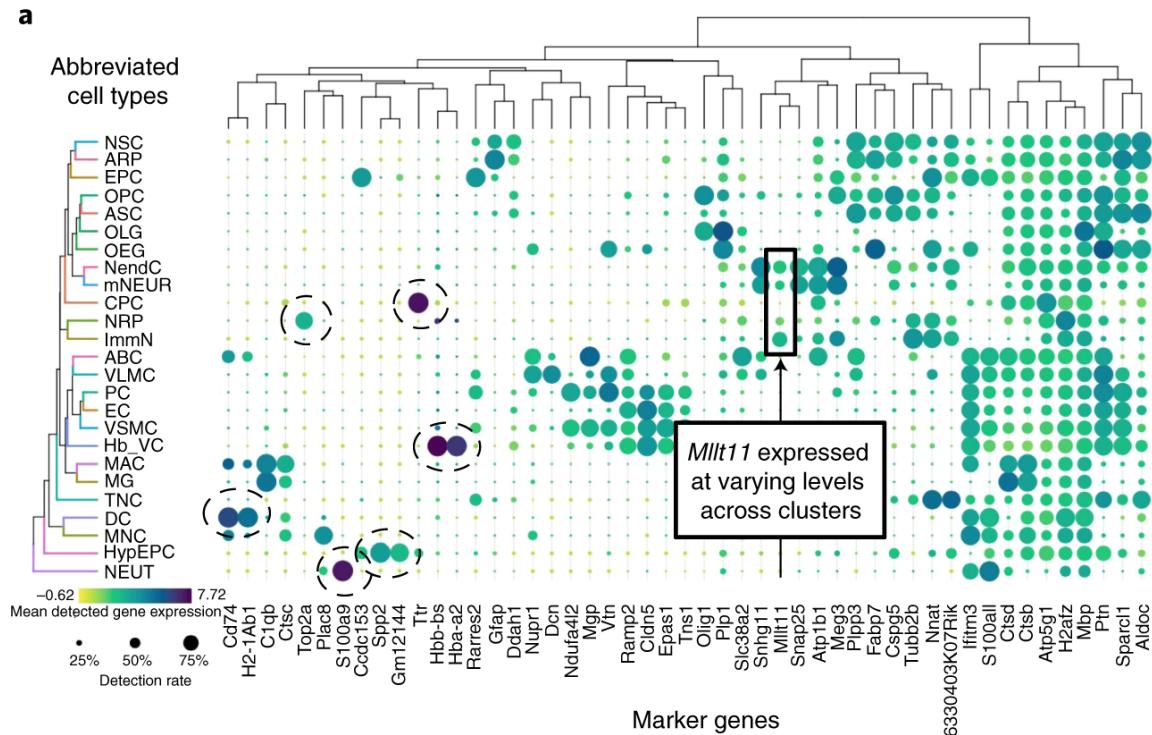
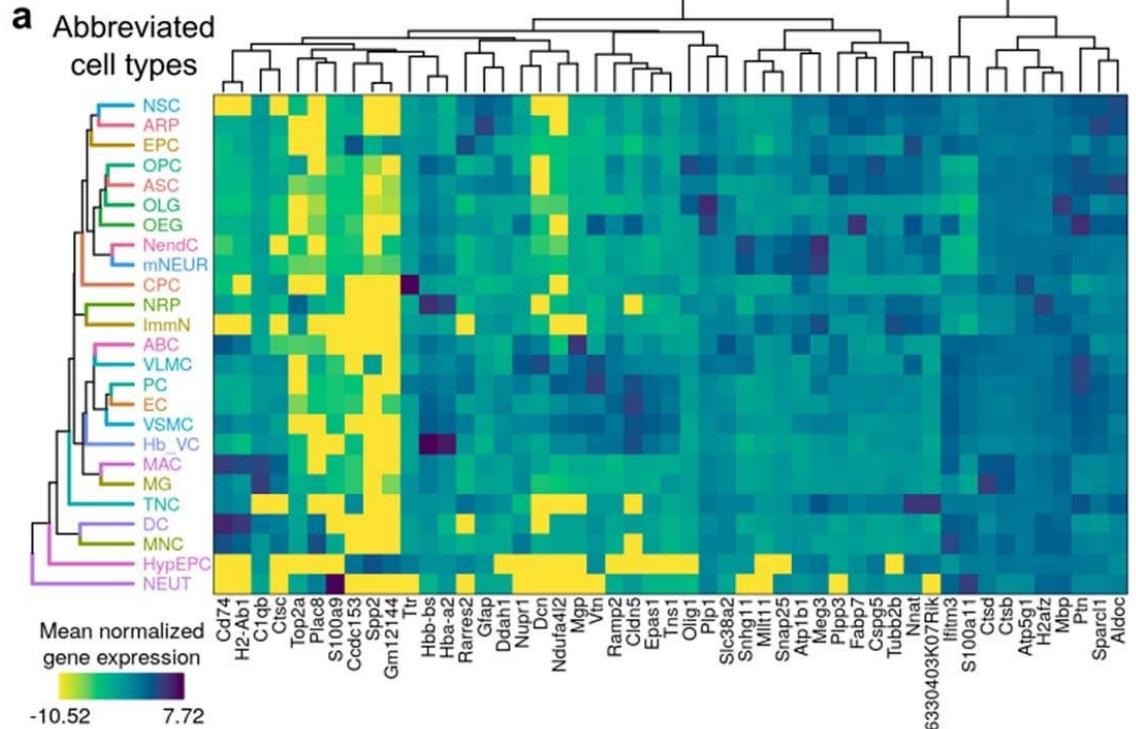
Easy



Challenging

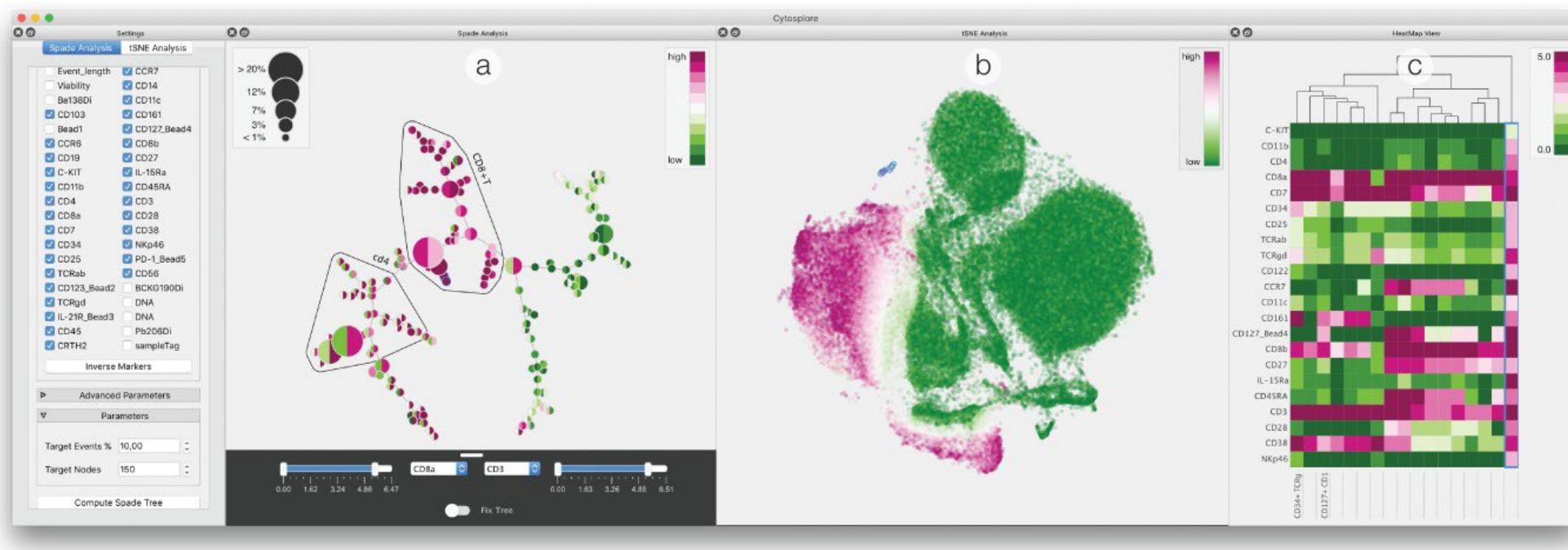


Alternatively: heatmaps & dot plots



Interactive visualization is important

- Interactive tools: Cytosplore, Loupe, cellxgene, ...
- Iterative visualization: Seurat, scanpy,...



Where do we get these marker genes?

Ideally: from a single cell atlas from a relevant organism, organ and disease context

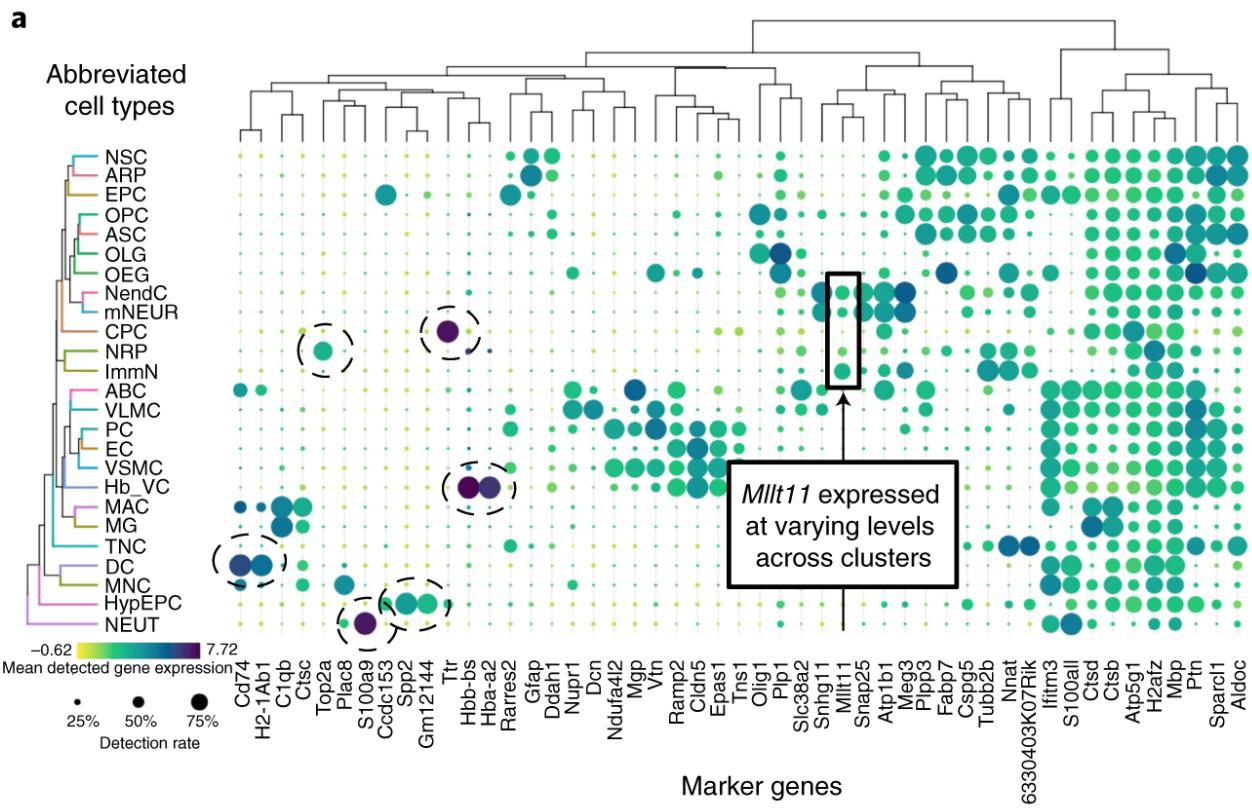
- “Expert knowledge”
- Literature
- Other scRNA-seq data
- Marker databases: PanglaoDB, CellMarker,...

Challenges:

- Few well-known markers
- Some well-known markers may not be as specific as expected

What if I don't have that many markers

- Identify “novel” markers by computing differential expression between a cluster and all other cells or between pairs of clusters
- Manually research differentially expressed genes to find functional information that may help identify the cell type



Complicating factors

1. Clusters that express markers of more than one cell type
 - Doublets?
 - Likely small, higher-than-average genes and UMIs per cell
 - Doublet detection tools: Scrublet, DoubletFinder, scds
2. Ambient RNA
 - RNA derived from one or more cell types that are sensitive to tissue dissociation
 - Markers of the contaminating cell types may be spread to all other cell clusters
 - Ambient RNA correction tools: SoupX, CellBlender

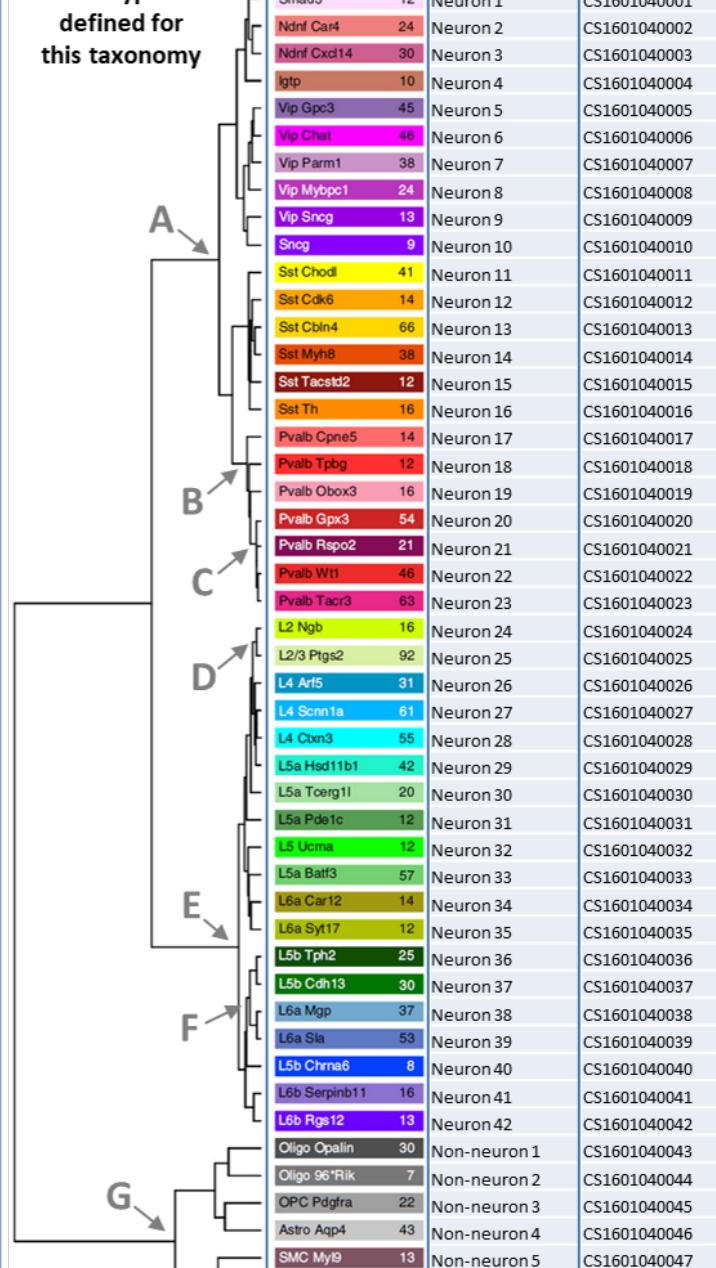
Watch out not to remove rare “interesting” cells!

Annotation verification

1. Using independent data (e.g. fluorescence in situ hybridization)
2. Multi-modal single-cell data
 - SNVs & CNVs
 - TCR/BCR
 - scRNA-seq+scATAC (mRNA + accessibility)
 - CITE-seq (surface proteins + mRNA)

Nomenclature

- How should we name cells?

Cell set nomenclature	Hierarchical organization of cell types defined for this taxonomy	Cell type alias	Cell type label	Cell type accession ID
A:		<i>Smad3</i> 12	Neuron 1	CS1601040001
B:		<i>Ndnf Ca4</i> 24	Neuron 2	CS1601040002
C:		<i>Ndnf Cxd14</i> 30	Neuron 3	CS1601040003
D:		<i>Igtp</i> 10	Neuron 4	CS1601040004
E:		<i>Vip Gpc3</i> 45	Neuron 5	CS1601040005
F:		<i>Vip Chat</i> 46	Neuron 6	CS1601040006
G:		<i>Vip Parm1</i> 38	Neuron 7	CS1601040007
		<i>Vip Mybpc1</i> 24	Neuron 8	CS1601040008
		<i>Vip Sneg</i> 13	Neuron 9	CS1601040009
		<i>Snog</i> 9	Neuron 10	CS1601040010
		<i>Sst Chodl</i> 41	Neuron 11	CS1601040011
		<i>Sst Cdk6</i> 14	Neuron 12	CS1601040012
		<i>Sst Cbln4</i> 66	Neuron 13	CS1601040013
		<i>Sst Myh8</i> 38	Neuron 14	CS1601040014
		<i>Sst Tacstd2</i> 12	Neuron 15	CS1601040015
		<i>Sst Th</i> 16	Neuron 16	CS1601040016
		<i>Pvalb Cpne5</i> 14	Neuron 17	CS1601040017
		<i>Pvalb Trpb</i> 12	Neuron 18	CS1601040018
		<i>Pvalb Obox3</i> 16	Neuron 19	CS1601040019
		<i>Pvalb Gpx3</i> 54	Neuron 20	CS1601040020
		<i>Pvalb Rspo2</i> 21	Neuron 21	CS1601040021
		<i>Pvalb Wt1</i> 46	Neuron 22	CS1601040022
		<i>Pvalb Tacr3</i> 63	Neuron 23	CS1601040023
		<i>L2 Ngb</i> 16	Neuron 24	CS1601040024
		<i>L2/3 Ptgs2</i> 92	Neuron 25	CS1601040025
		<i>L4 Arf5</i> 31	Neuron 26	CS1601040026
		<i>L4 Scnn1a</i> 61	Neuron 27	CS1601040027
		<i>L4 Ctxn3</i> 55	Neuron 28	CS1601040028
		<i>L5a Hsd11b1</i> 42	Neuron 29	CS1601040029
		<i>L5a Tcerg11</i> 20	Neuron 30	CS1601040030
		<i>L5a Pde1c</i> 12	Neuron 31	CS1601040031
		<i>L5 Ucma</i> 12	Neuron 32	CS1601040032
		<i>L5a Batt3</i> 57	Neuron 33	CS1601040033
		<i>L6a Car12</i> 14	Neuron 34	CS1601040034
		<i>L6a Syt17</i> 12	Neuron 35	CS1601040035
		<i>L5b Tph2</i> 25	Neuron 36	CS1601040036
		<i>L5b Cd13</i> 30	Neuron 37	CS1601040037
		<i>L6a Mgp</i> 37	Neuron 38	CS1601040038
		<i>L6a Sla</i> 53	Neuron 39	CS1601040039
		<i>L5b Chrna6</i> 8	Neuron 40	CS1601040040
		<i>L6b Serpinb11</i> 16	Neuron 41	CS1601040041
		<i>L6b Rgs12</i> 13	Neuron 42	CS1601040042
		<i>Oligo Opalin</i> 30	Non-neuron 1	CS1601040043
		<i>Oligo 96'Rik</i> 7	Non-neuron 2	CS1601040044
		<i>OPC Pdgfra</i> 22	Non-neuron 3	CS1601040045
		<i>Astro Aqp4</i> 43	Non-neuron 4	CS1601040046
		<i>SMC My9</i> 13	Non-neuron 5	CS1601040047
		<i>Endo Xdh</i> 14	Non-neuron 6	CS1601040048
		<i>Micro Ctss</i> 22	Non-neuron 7	CS1601040049

Summary

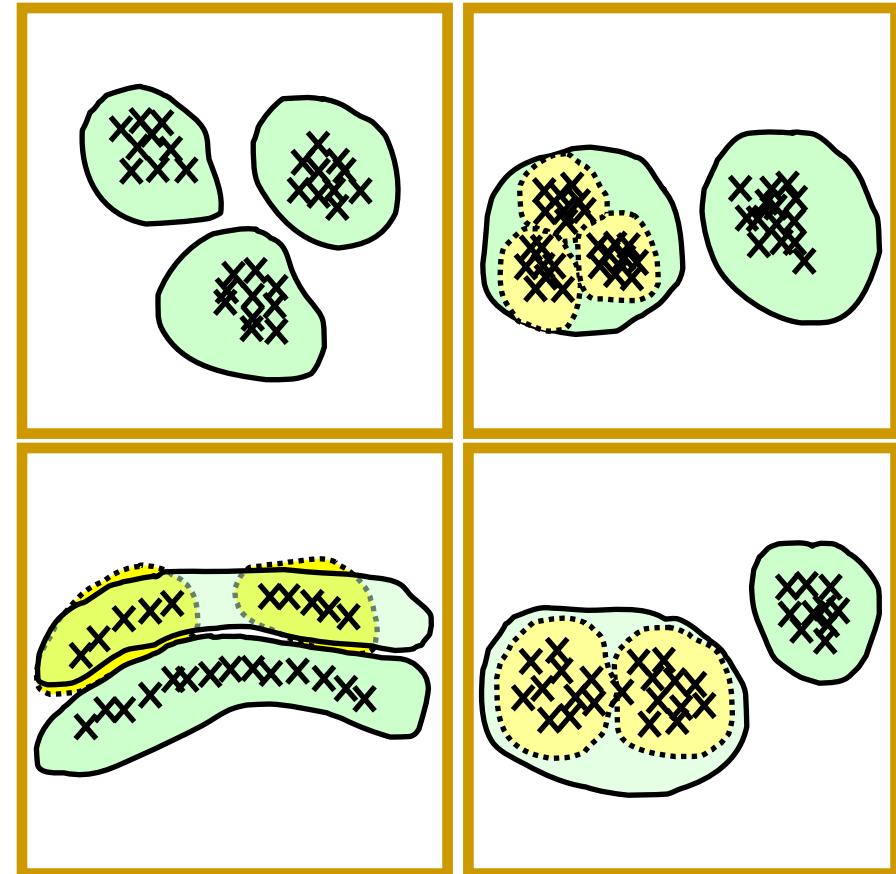
- Start by identifying major well-known cell types (clearly defined, discrete cell clusters)
- Split the data into broad subsets (e.g., immune, endothelial and tumor) and analyze each separately
- Cell subtypes or poorly defined clusters are challenging
- Manual annotations heavily rely on marker genes and expert knowledge

Outline

- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- scRNA-seq clustering
- Annotating clusters
- Cluster validation
- Challenges & outlook

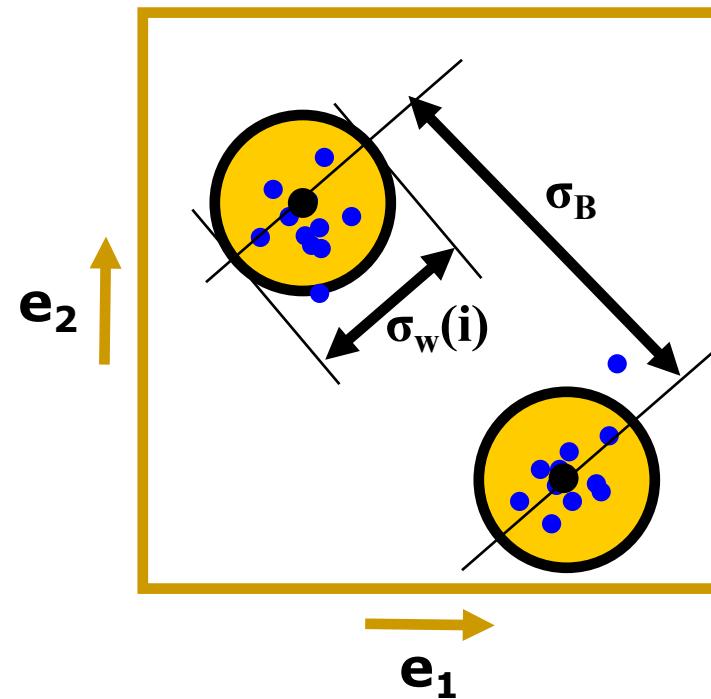
Clustering is subjective!

- Principle choices
 - Similarity measure
 - Algorithm
- Different choice leads to different results
 - Subjectivity becomes reality
- Cluster process
 - Validate, interpret (generate hypothesis), repeat steps



Cluster criteria

- Silhouette score
 - Goal: optimize cohesion within a cluster and separation between clusters
 - Seek: clustering that maximizes SI



Silhouette score

1. Mean distance between i and all other points in cluster C_i

$$a(i) = \frac{1}{|C_i| - 1} \sum_{j \in C_i, i \neq j} d(i, j)$$

2. Mean nearest cluster distance of i

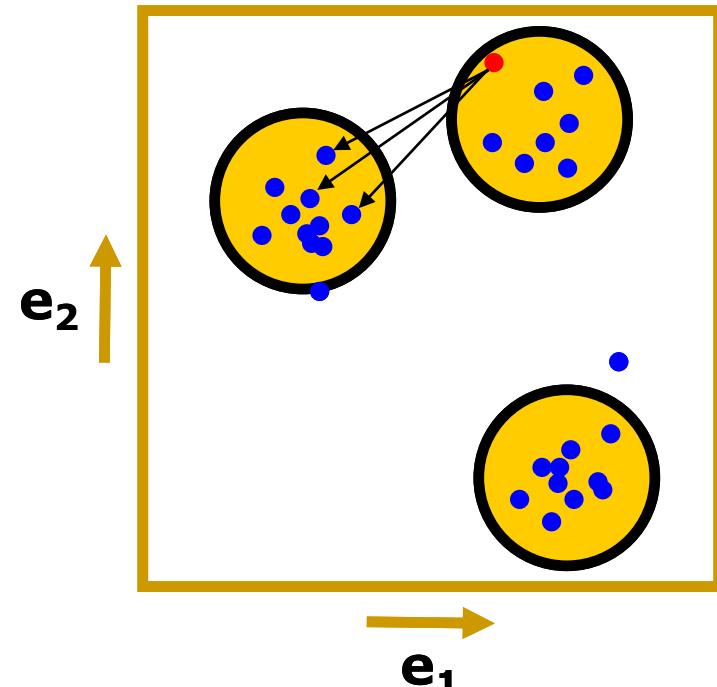
$$b(i) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j)$$

3. Silhouette score for i

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

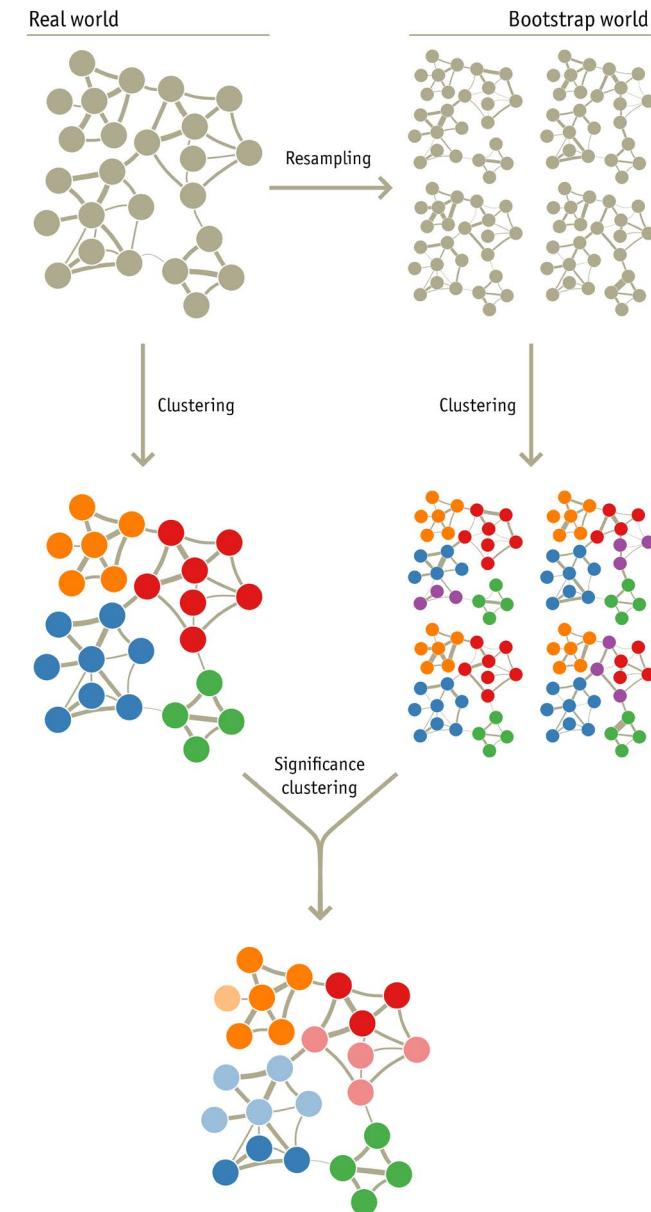
4. Total silhouette score

$$SI = \frac{1}{N} \sum s(i)$$



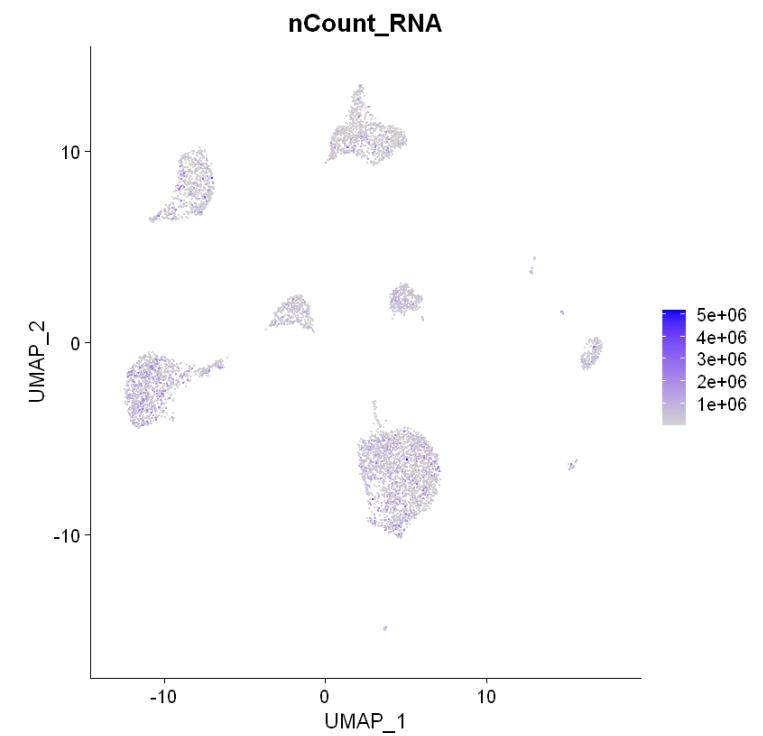
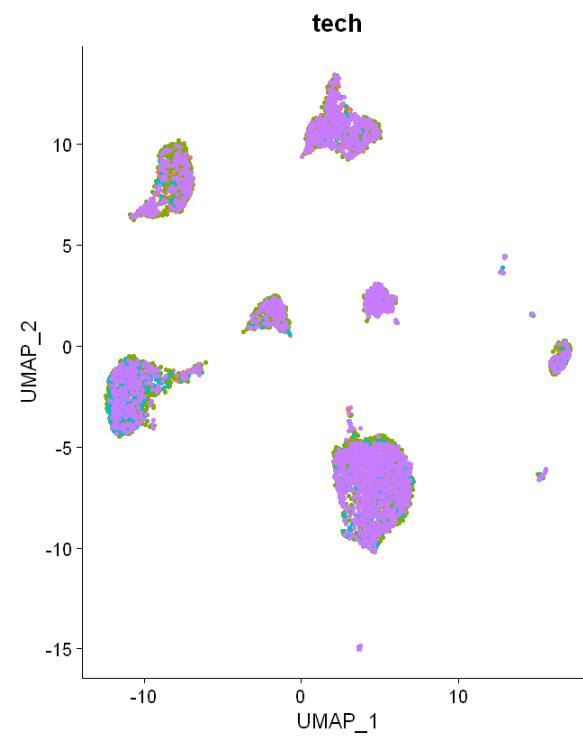
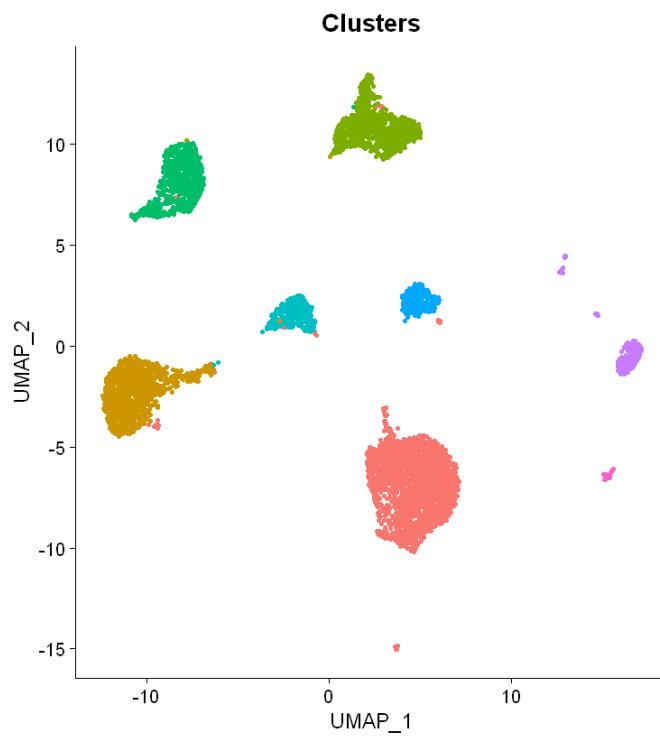
Bootstrapping

- How confident can you be that the clusters you see are real?
 - Take a random set of cells
 - Cluster
 - Compare to original clustering
 - Estimate support for clustering



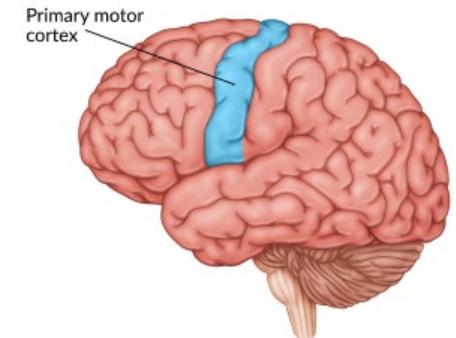
Always check QC data

- Are your clusters mainly related to batches, qc-measures (especially detected genes)?

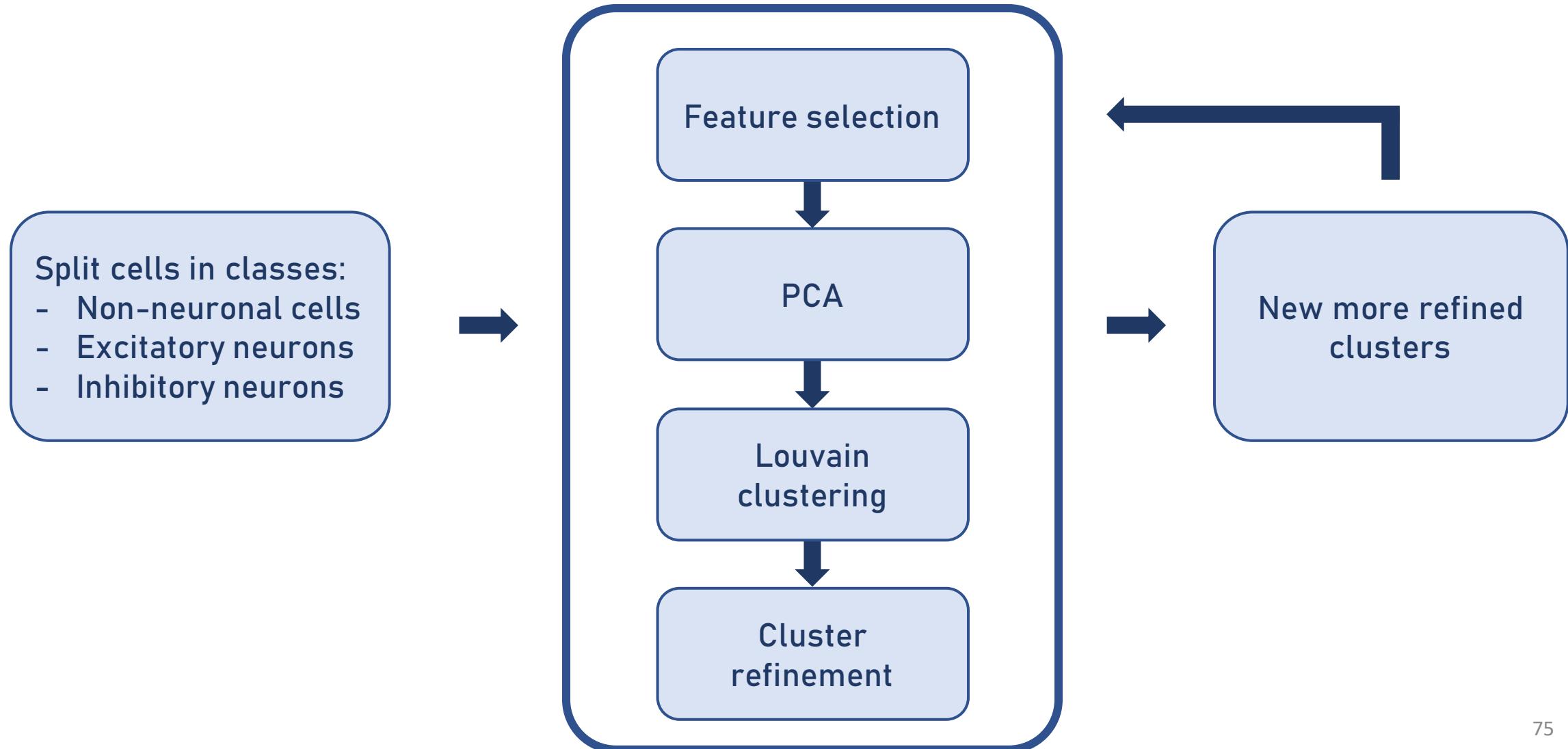


Example: annotating human brain cells

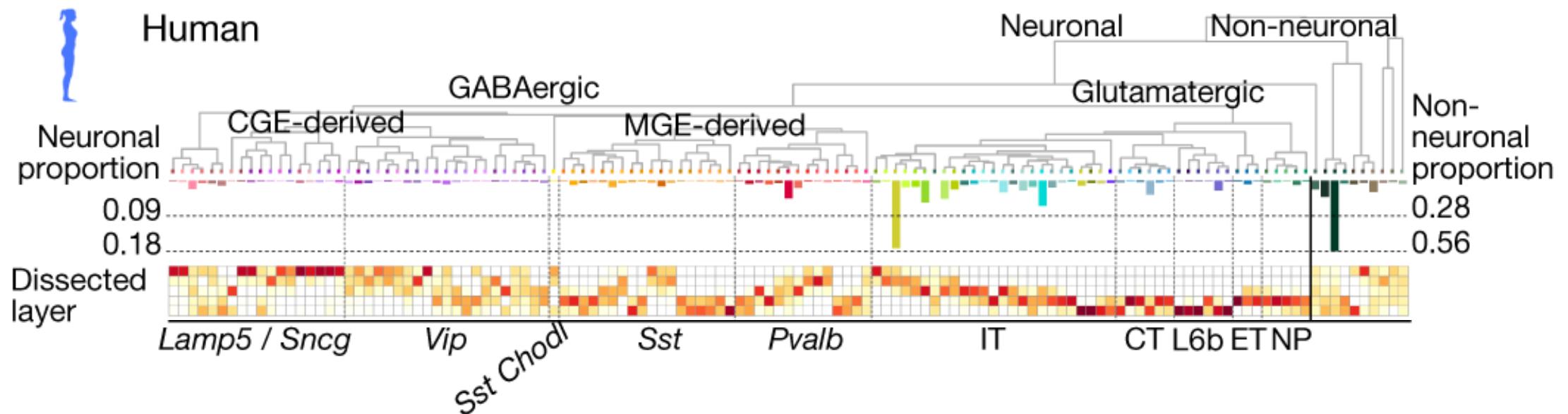
72,621 cells
32,991 genes
127 clusters



Iterative clustering approach



Example: annotating human brain cells



Outline

- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- scRNA-seq clustering
- Annotating clusters
- Cluster validation
- Challenges & outlook

Challenges

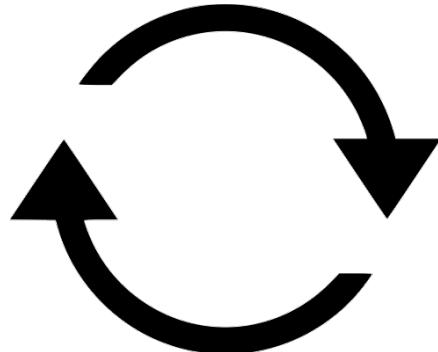
- **Subjectivity:** what is a cell type?
 - Different parameters yield different results
 - Validation is important
- **Scalability:** number of cells has grown from $\sim 10^2$ to $\sim 10^6$
 - Computational efficiency
 - Visual exploration, crowding problem

Downside of clustering

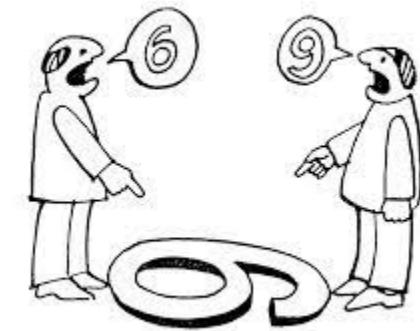
Time consuming



Not reproducible

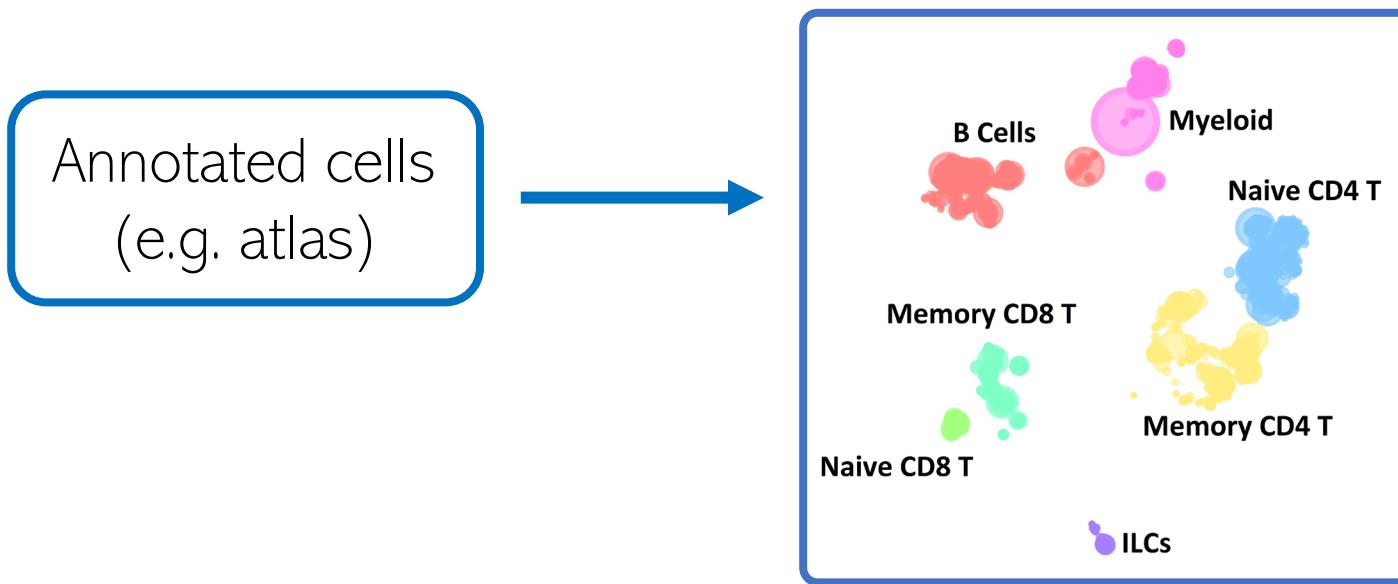


Subjective

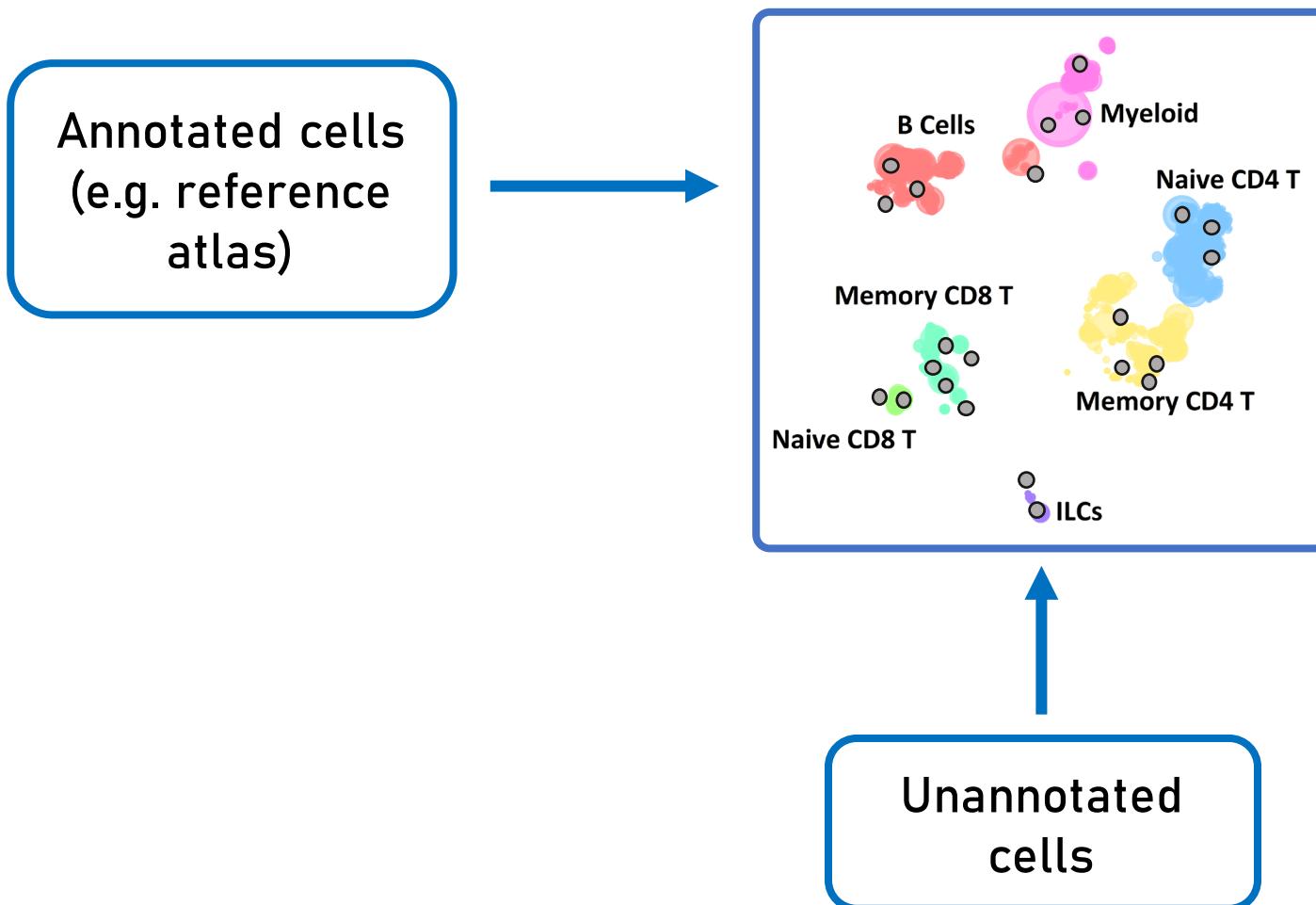


Lots of single-cell data is available nowadays!
Can we use that to annotate our cells?

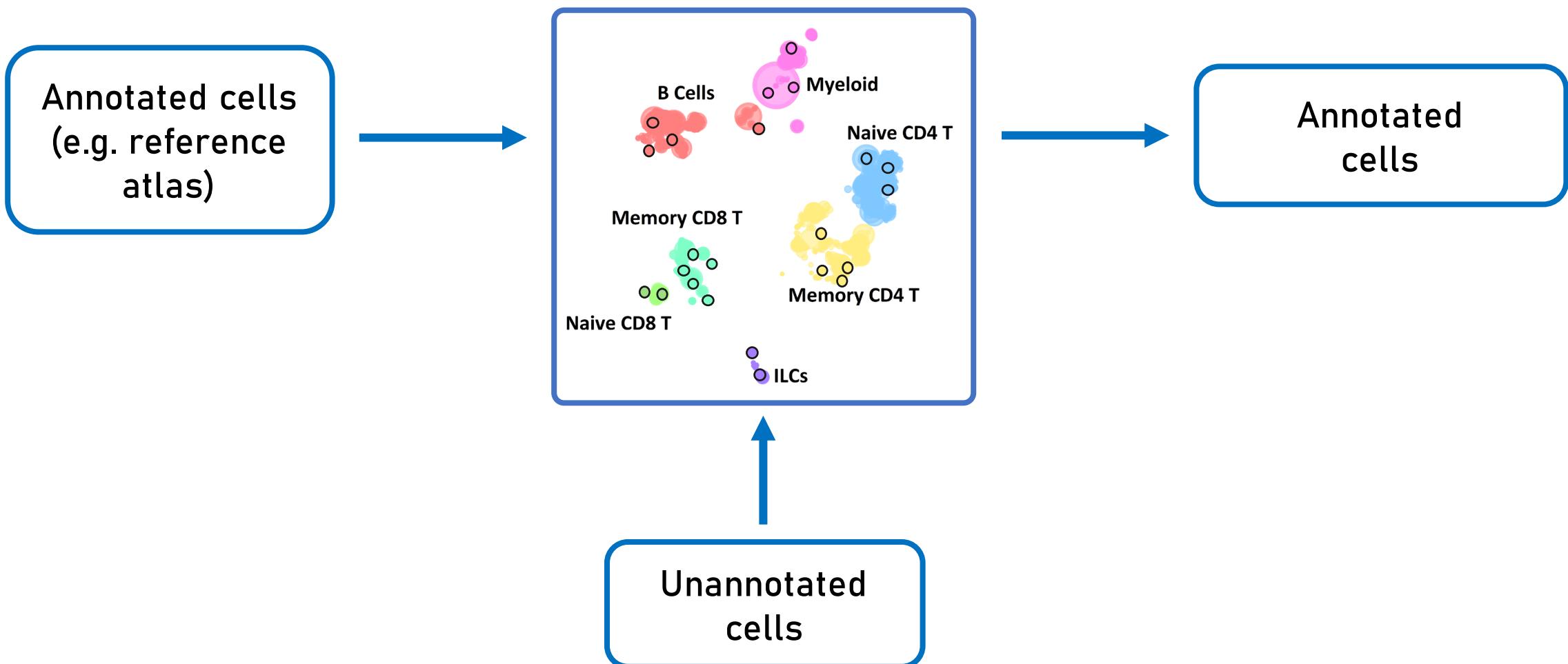
Supervised approach



Supervised approach



Supervised approach



Clustering practical

- Hierarchical clustering: distances and linkage methods
- k -Means
- Graph-based clustering
- Annotating clusters

Resources

- Kiselev et al. "Challenges in unsupervised clustering of single-cell RNA-seq data"
<https://doi.org/10.1038/s41576-018-0088-9>
- Duò et al. "A systematic performance evaluation of clustering methods for single-cell RNA-seq data"
<https://doi.org/10.12688/f1000research.15666.2>
- Orchestrating Single-Cell Analysis with Bioconductor
<https://osca.bioconductor.org/>
- Hemberg single cell course: Analysis of single cell RNA-seq data
<https://scrnaseq-course.cog.sanger.ac.uk/website/index.html>
- Slides Åsa Björklund (NBIS, SciLifeLab)
<https://github.com/NBISweden/workshop-scRNASeq/tree/master/slides2019>
- Tutorial: guidelines for annotating single-cell transcriptomic maps using automated and manual methods
<https://doi.org/10.1038/s41596-021-00534-0>