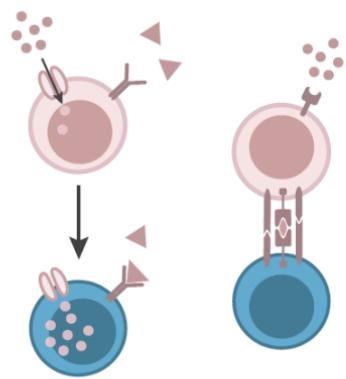


Single Cell RNA-seq Clustering

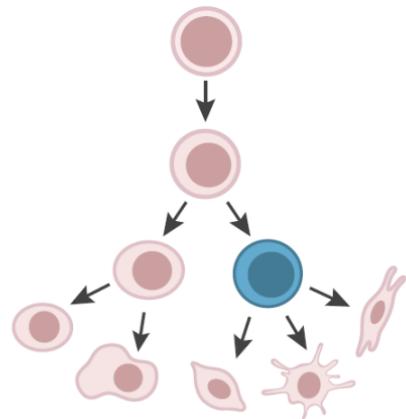
Marcel Reinders &
Lieke Michielsen & Ahmed Mahfouz
Delft Bioinformatics Lab, TU Delft
Leiden Computational Biology Center, LUMC

Cell Identity determined by diverse factors

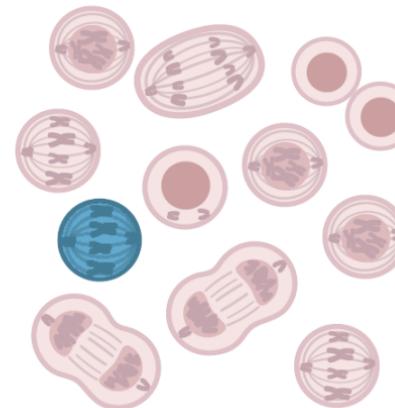
Environmental stimuli



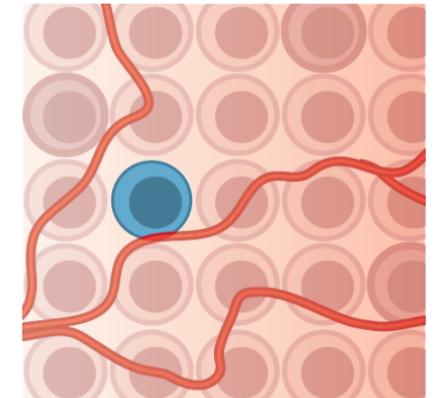
Cell development



Cell cycle

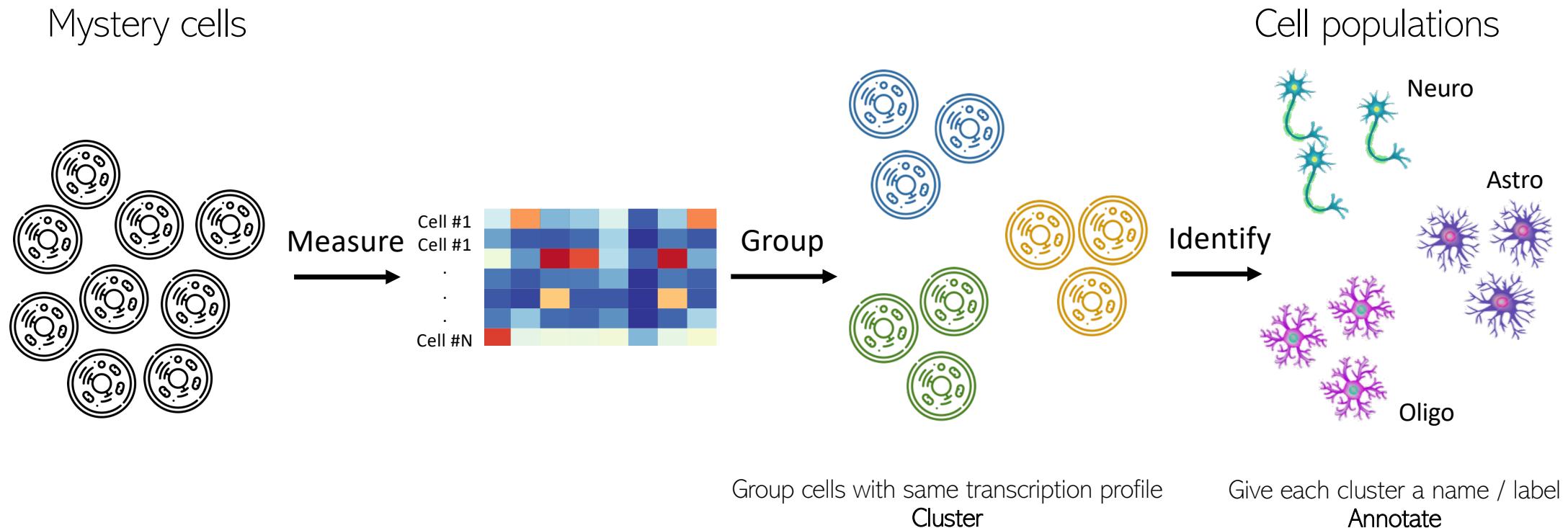


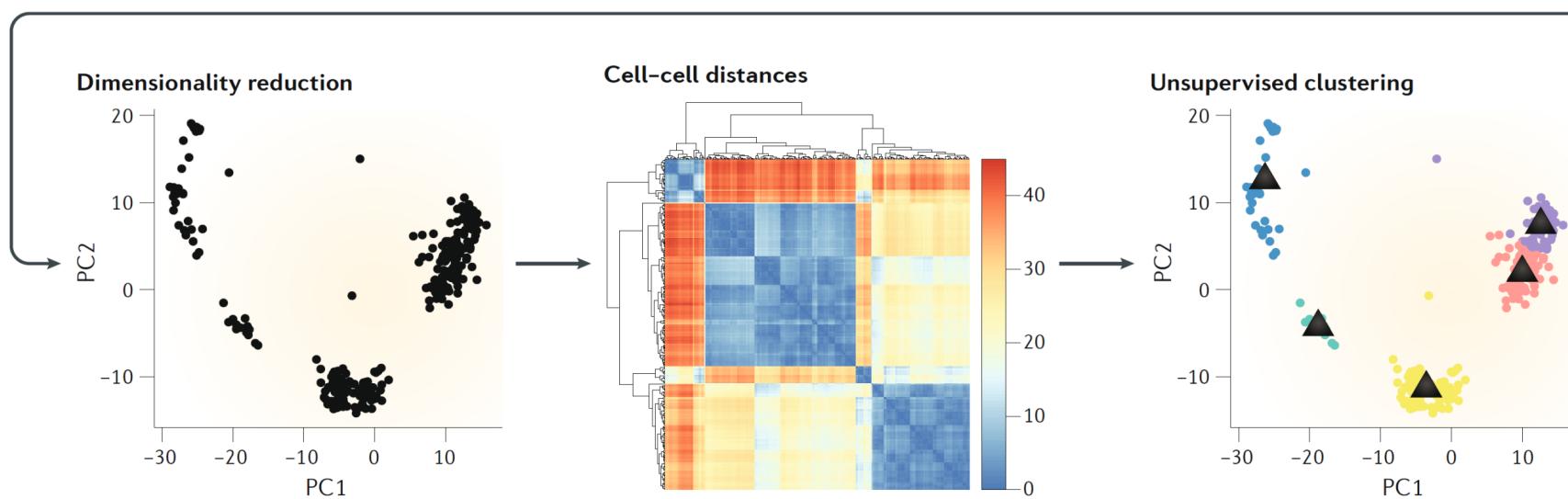
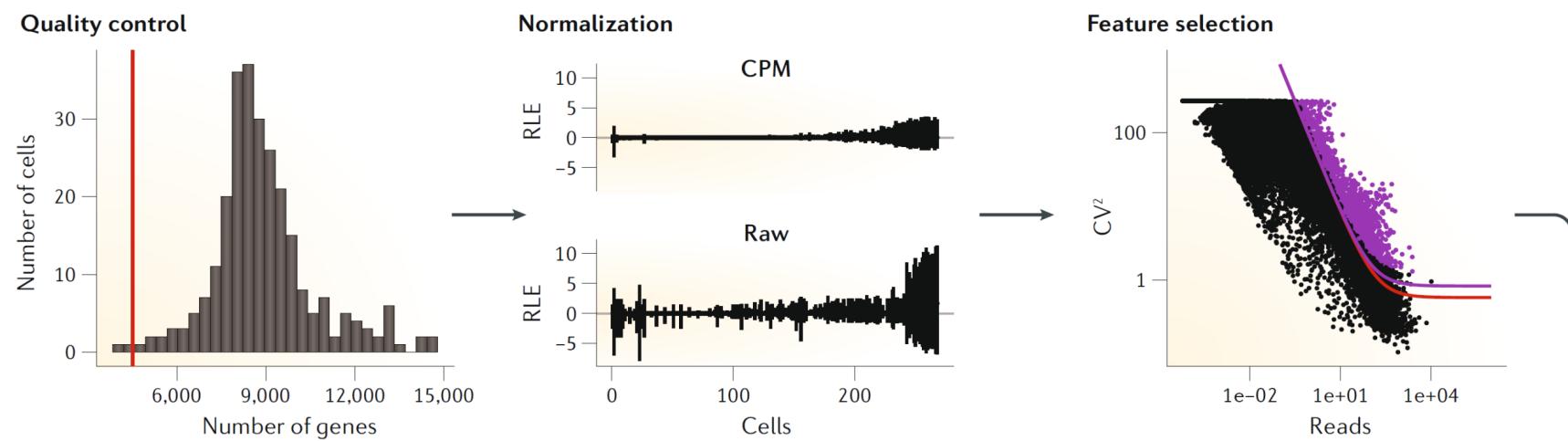
Spatial context

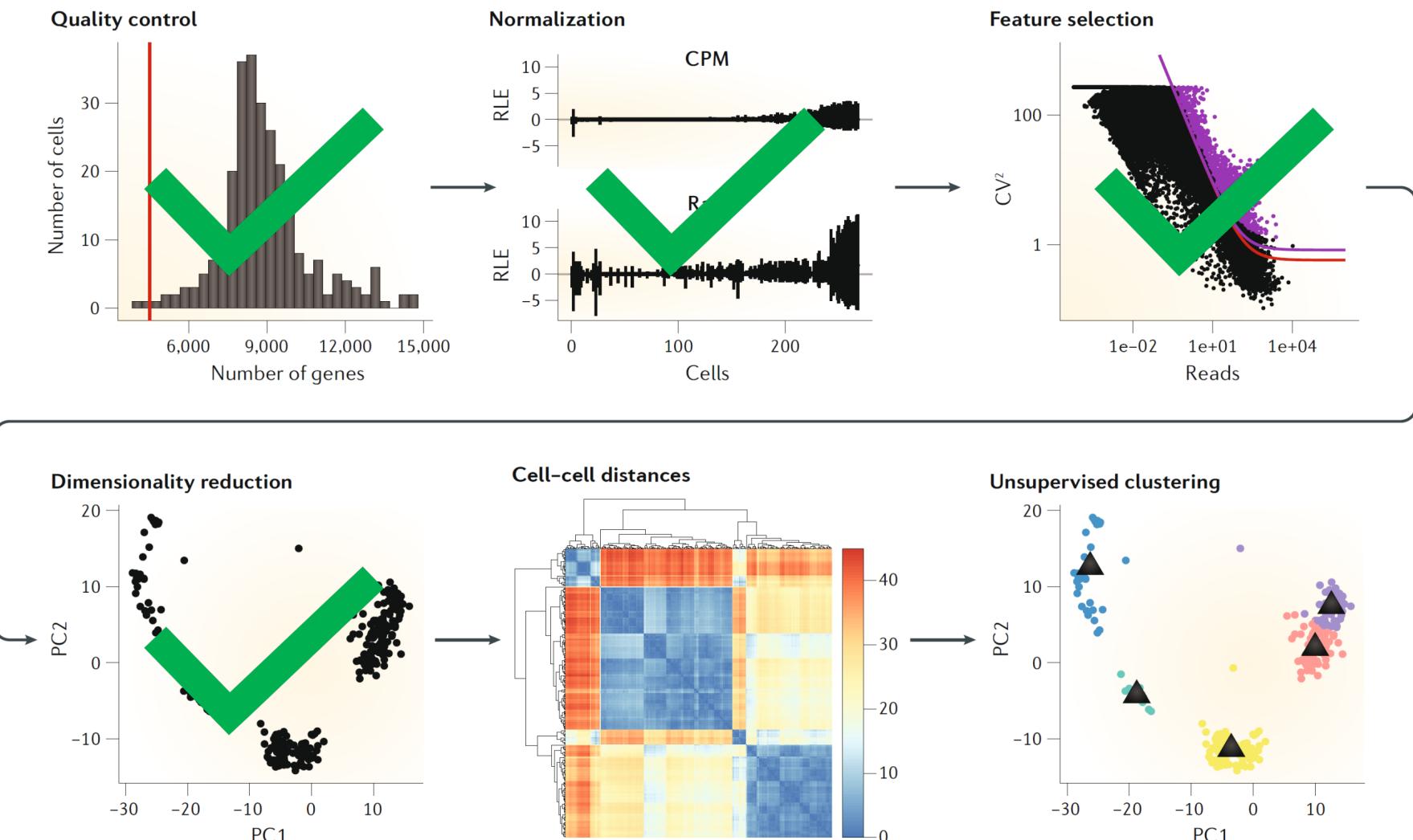


- **Cell type:** Permanent state (e.g., a hepatocyte typically cannot turn into a neuron).
- **Cell state:** Transiently state (e.g., during differentiation, stimulation, spatial patterning)

How can we identify cell populations?







Outline

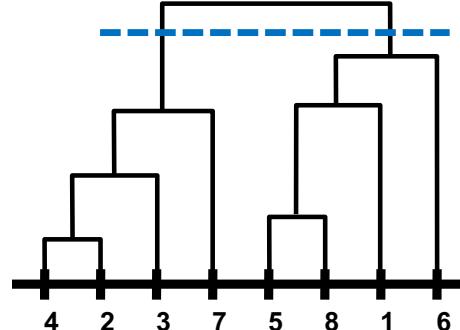
- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- Cluster validation
- scRNA-seq clustering
- Annotating clusters

Outline

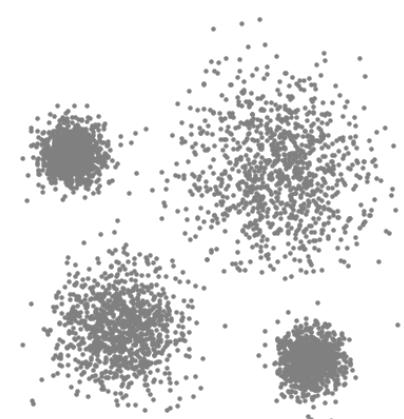
- **Introduction to clustering**

- Hierarchical clustering
- k -Means clustering
- Graph-based clustering
- Cluster validation
- scRNA-seq clustering
- Annotating clusters

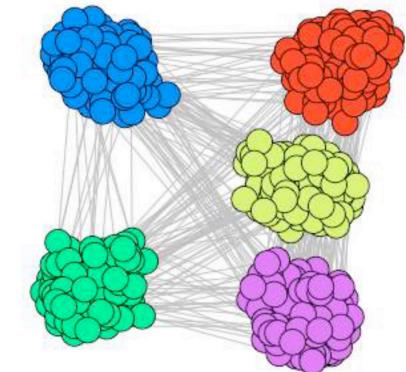
Many clustering approaches



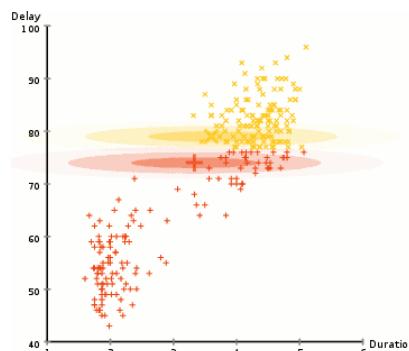
Hierarchical clustering



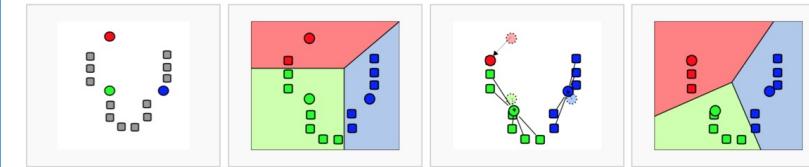
Mean shift clustering



Graph-based clustering



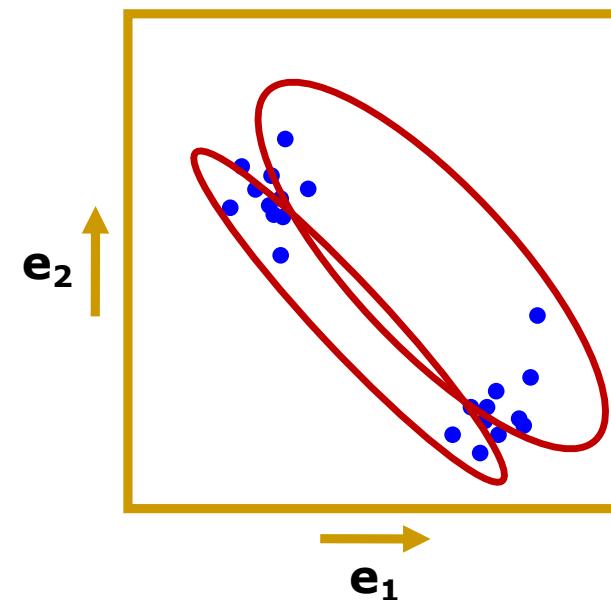
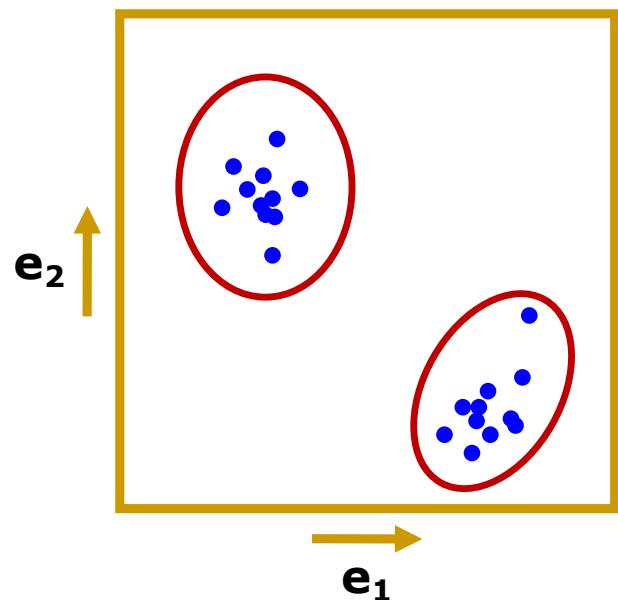
Gaussian mixture modeling



k -means clustering

Clustering

What defines a good clustering



Clustering

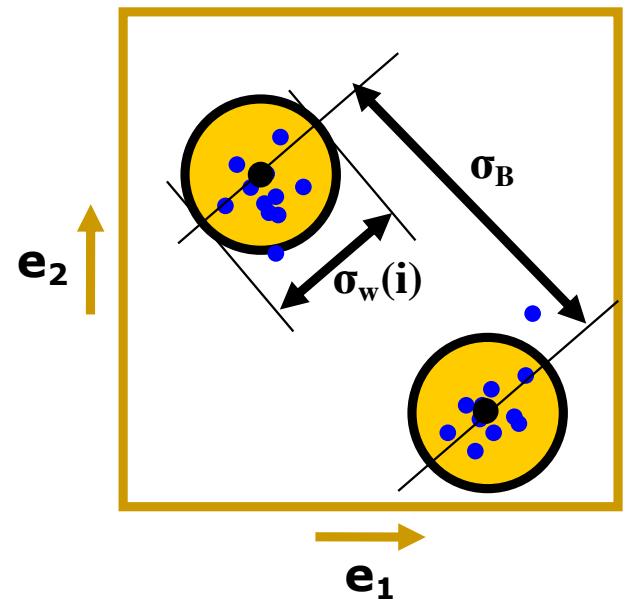
What defines a good clustering

Structure when:

1. Samples within cluster resemble each other
*(small pairwise distance between cells within cluster:
small within variance, $\sigma_w(i)$, for each cluster i)*
2. Clusters deviate from each other
*(large pairwise distance between cells of different
clusters: large between variance, σ_B)*

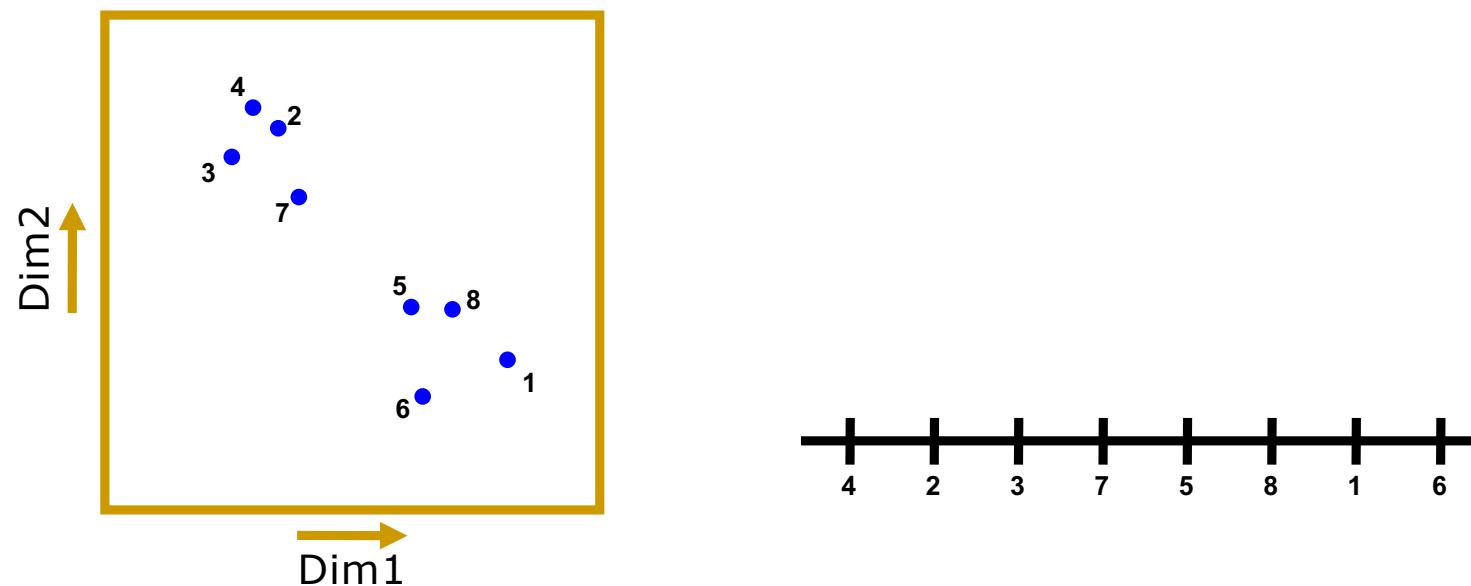
Group sa

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



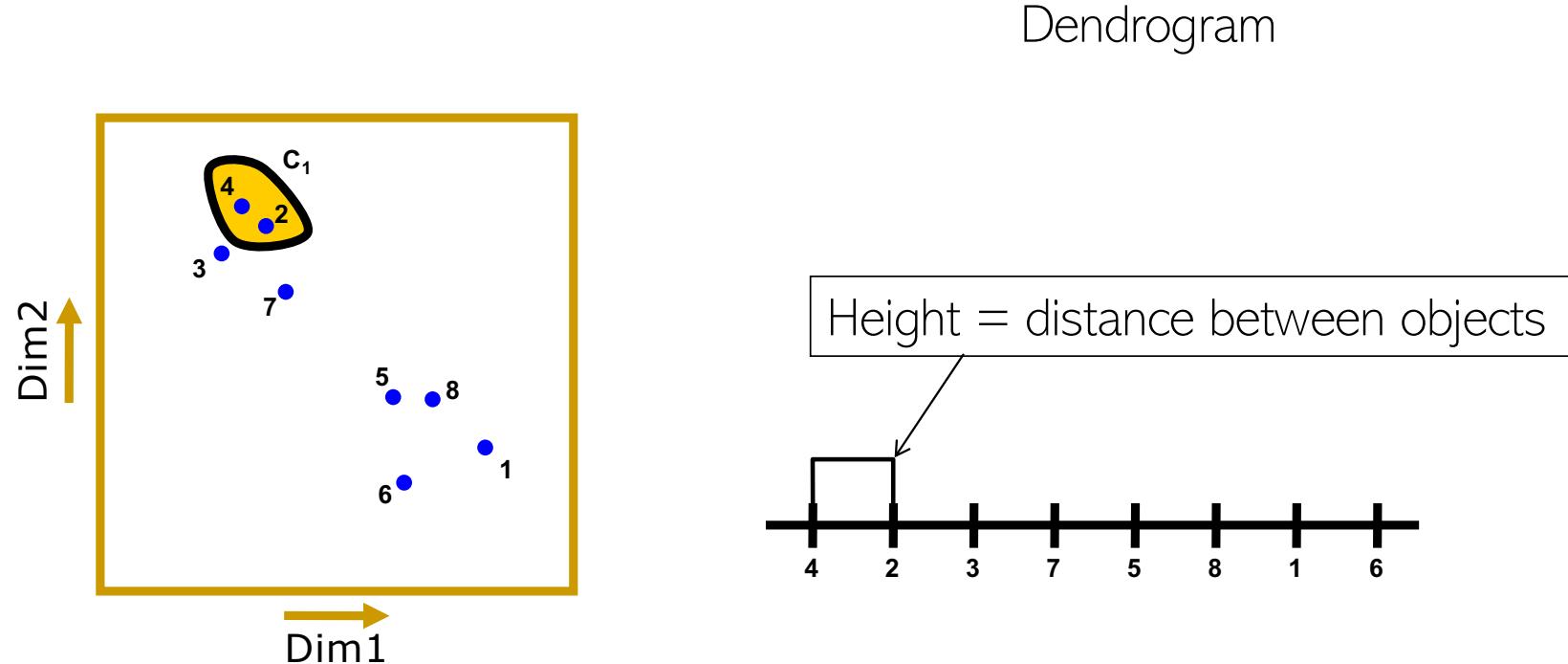
Hierarchical clustering

Hierarchical clustering



Find most similar objects (cells) and group them

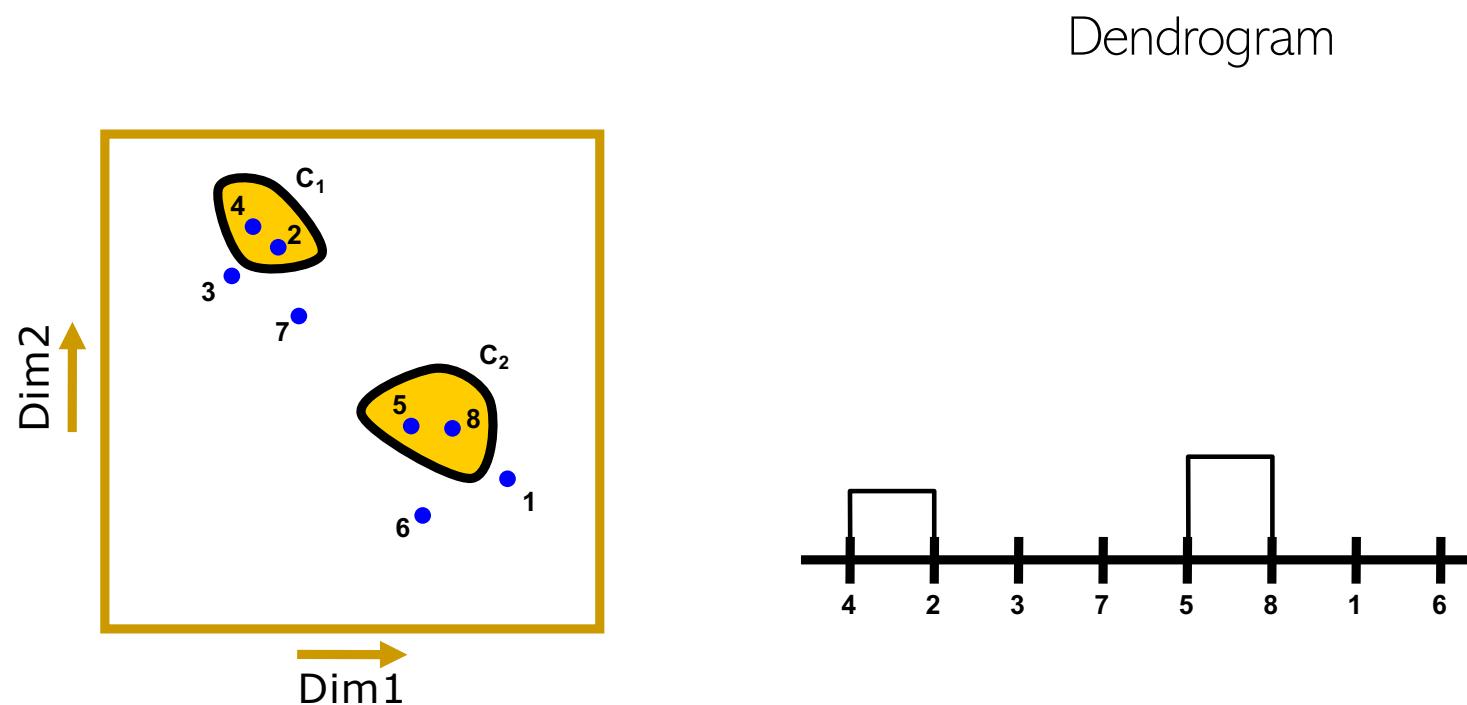
Hierarchical clustering



These are: objects 4 and 2

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

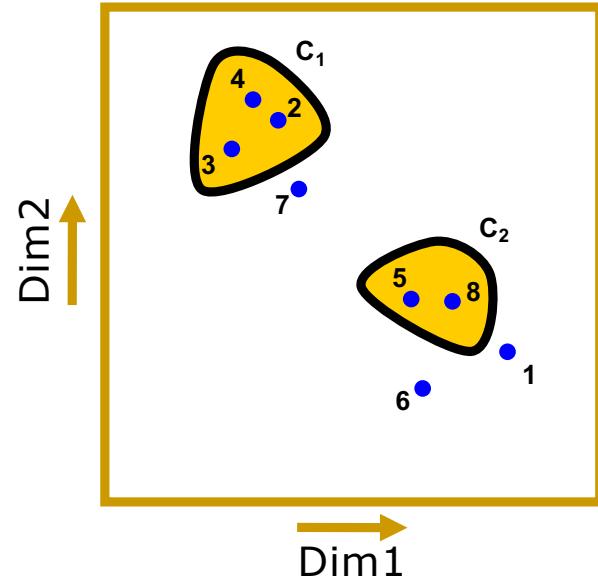
Hierarchical clustering



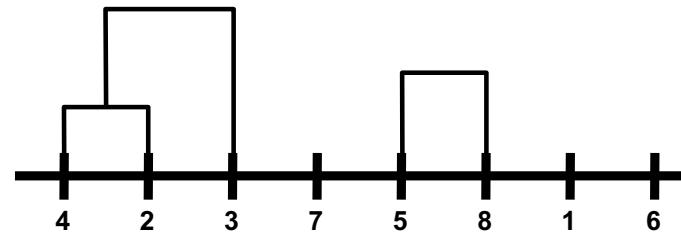
These are: objects 5 and 8

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

Hierarchical clustering

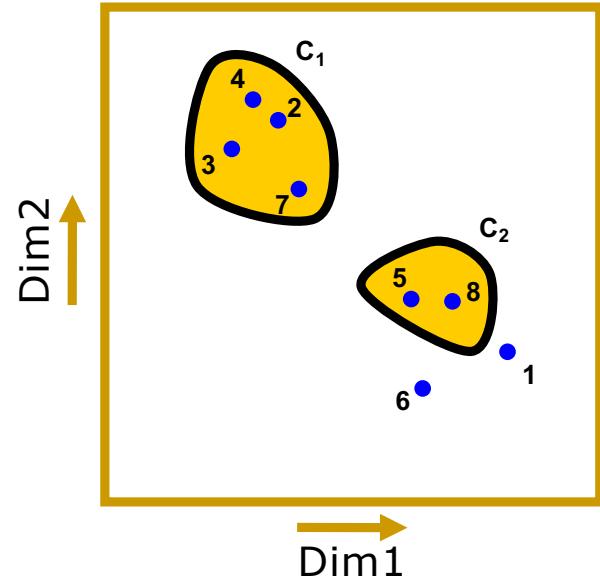


Dendrogram

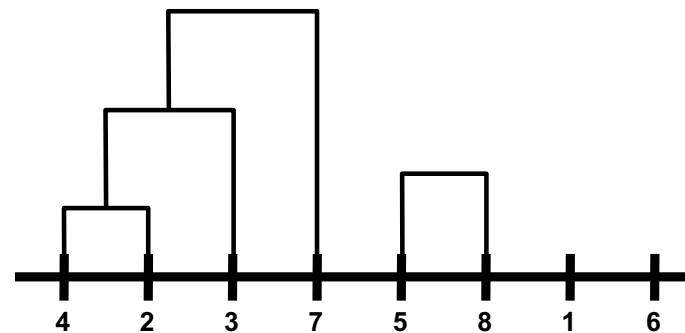


Join object 3 and cluster 1
Repeat process

Hierarchical clustering

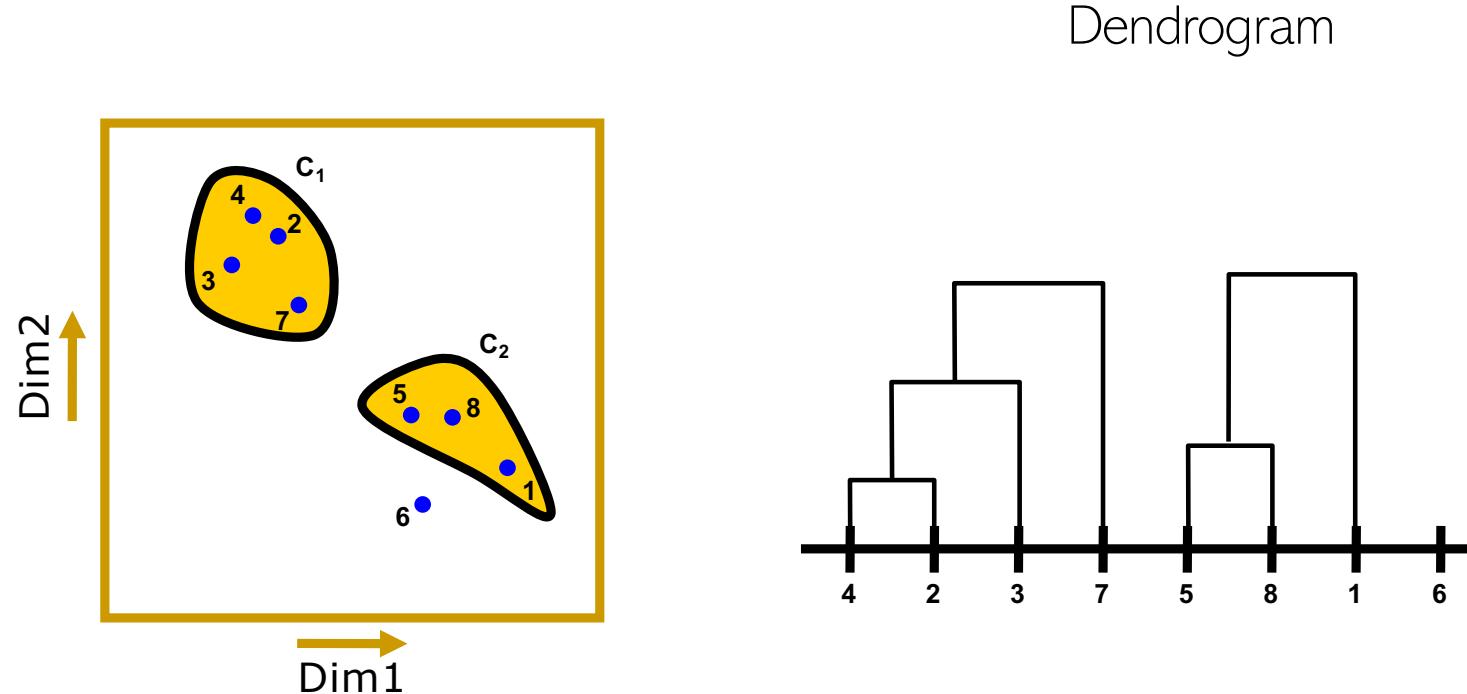


Dendrogram



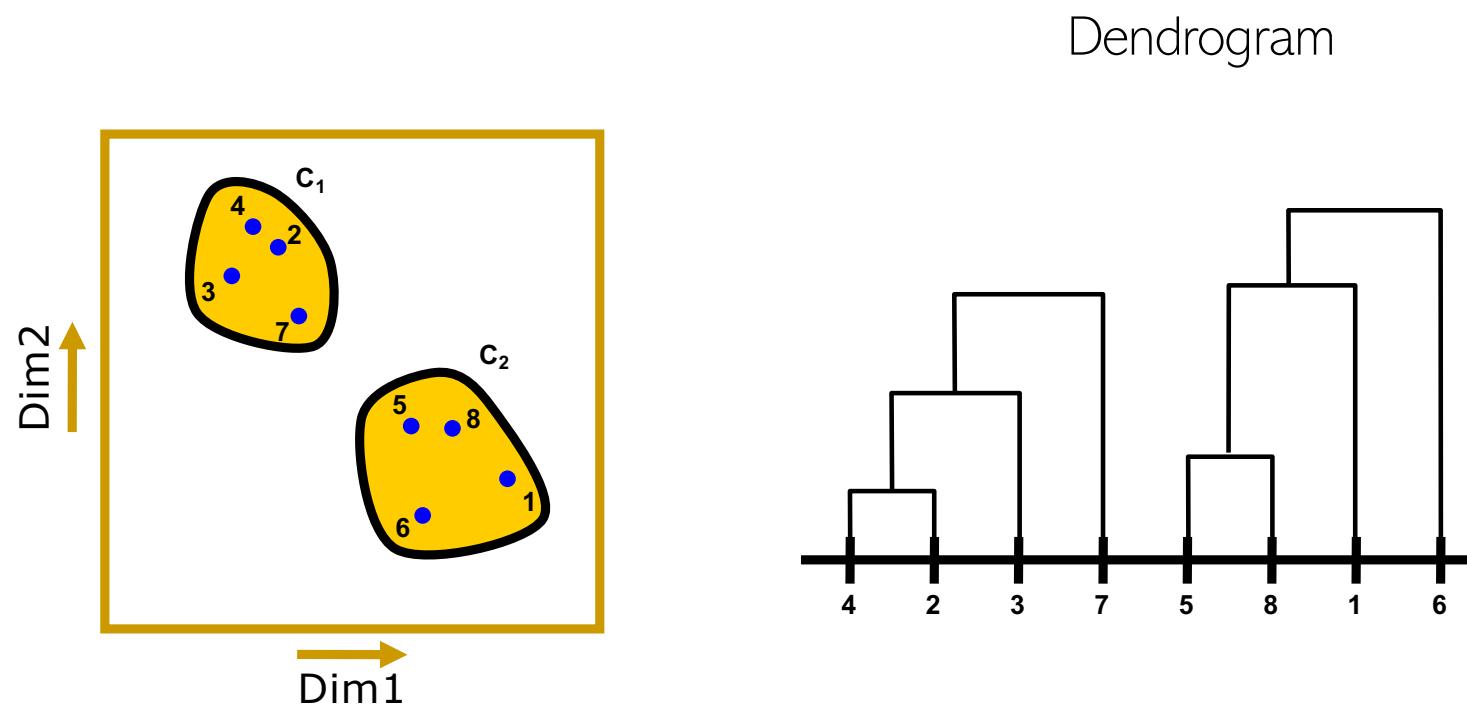
Join [object 7 and cluster 1] -> [cluster 1]
Repeat process

Hierarchical clustering



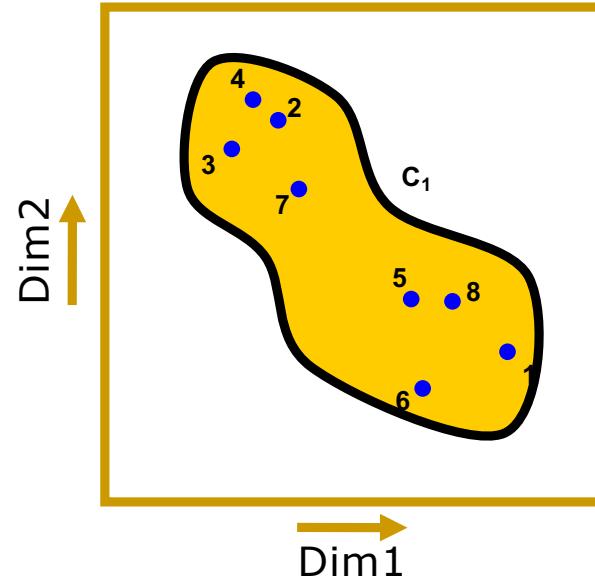
Join [object 1 and cluster 2] -> [cluster 2]
Repeat process

Hierarchical clustering

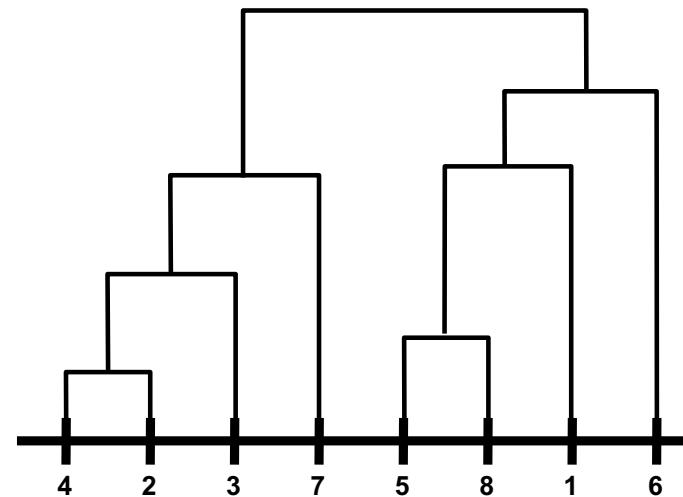


Join [object 6 and cluster 2] -> [cluster 2]
Repeat process

Hierarchical clustering

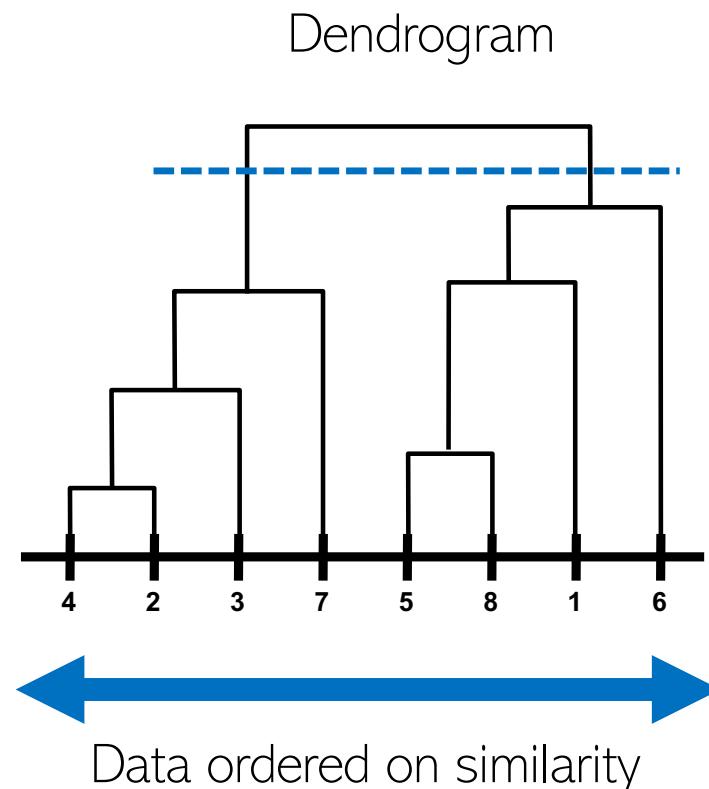
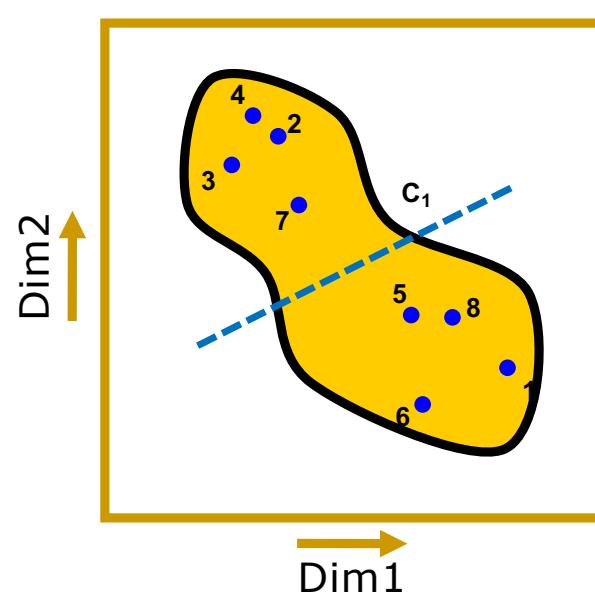


Dendrogram



Join [cluster 1 and cluster 2] -> [cluster 1]
All in one cluster: FINISHED!

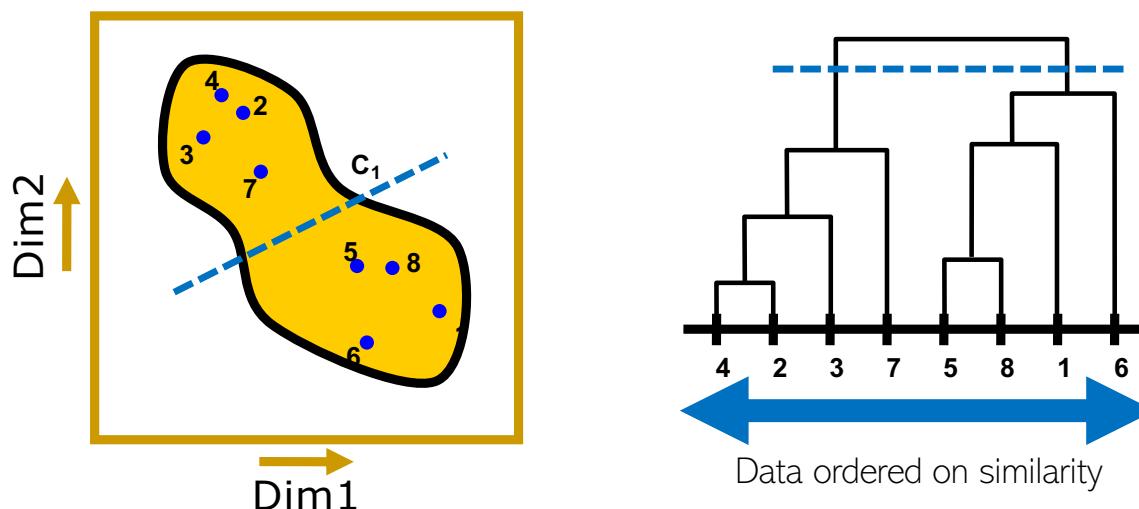
Hierarchical clustering



Hierarchical clustering

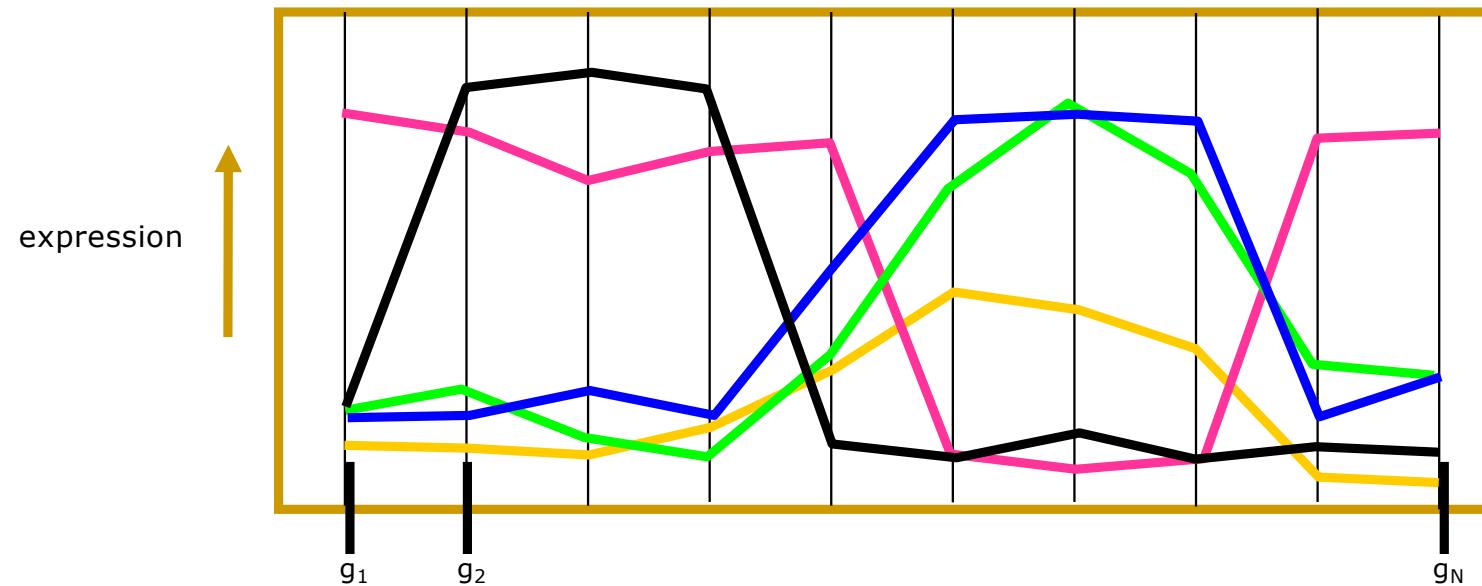
Need to know:

- Similarity between objects
- Similarity between clusters



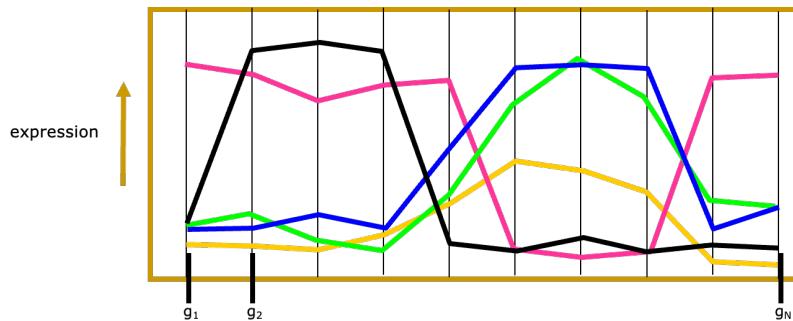
Hierarchical clustering

Similarity between objects



Hierarchical clustering

Similarity between objects



Euclidean distance

$$d(g_i, g_j) = \sqrt{(\sum ((x_i - x_j)^2))}$$

$$\begin{aligned} d(\bullet, \bullet) &< d(\bullet, \circ) \\ d(\bullet, \bullet) &<< d(\bullet, \circ) \\ d(\bullet, \bullet) &<< d(\bullet, \bullet) \end{aligned}$$

Match exact shape

Pearson correlation

$$1 - \rho_{ij}$$

$$\begin{aligned} d(\bullet, \bullet) &\approx d(\bullet, \circ) \\ d(\bullet, \bullet) &<< d(\bullet, \circ) \\ d(\bullet, \bullet) &<< d(\bullet, \bullet) \end{aligned}$$

Ignore amplitude

Mixed Pearson correlation

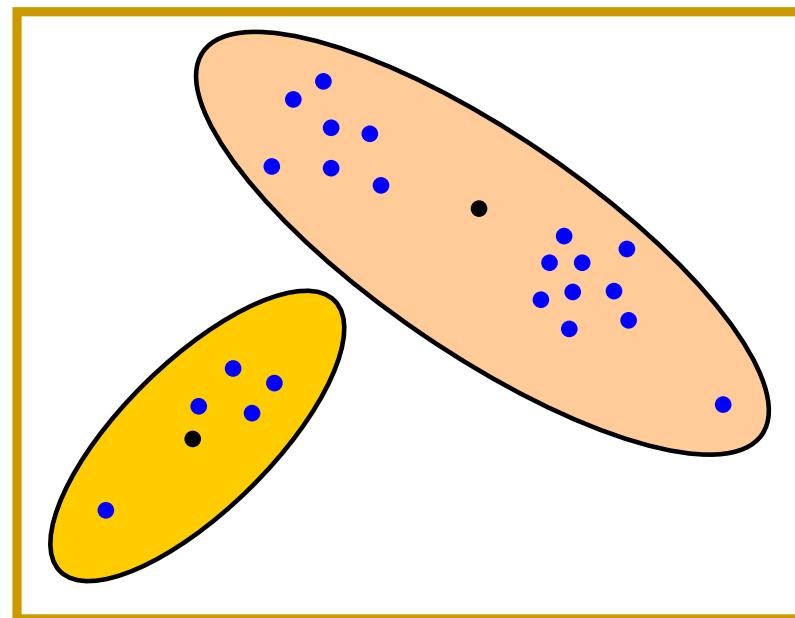
$$1 - |\rho_{ij}|$$

$$\begin{aligned} d(\bullet, \bullet) &\approx d(\bullet, \circ) \\ d(\bullet, \bullet) &\approx d(\bullet, \circ) \\ d(\bullet, \bullet) &<< d(\bullet, \bullet) \end{aligned}$$

Ignore amplitude and sign

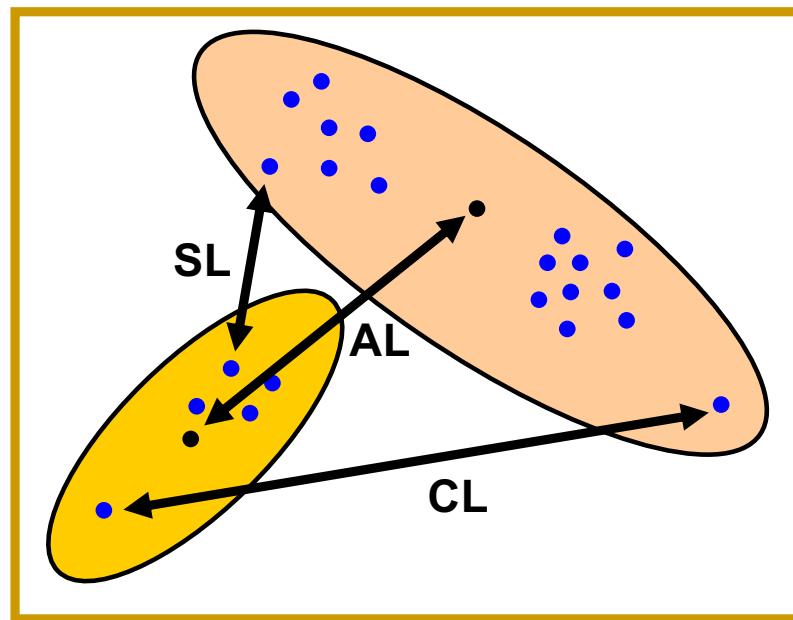
Hierarchical clustering

Similarity between clusters



Hierarchical clustering

Similarity between clusters



Single linkage:

Closest objects

Complete linkage:

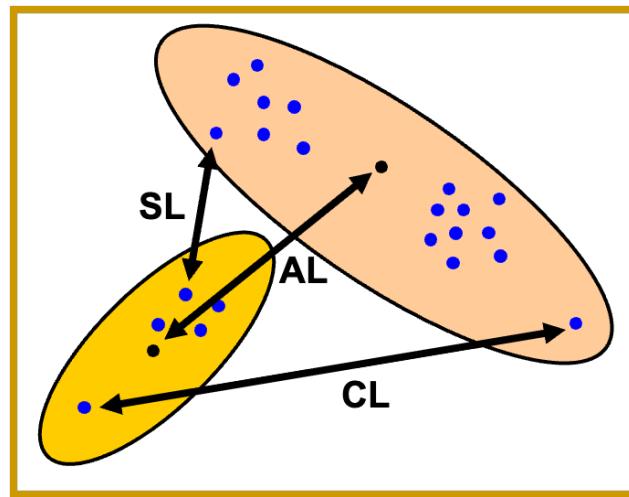
Furthest objects

Average linkage:

Average similarity

Hierarchical clustering

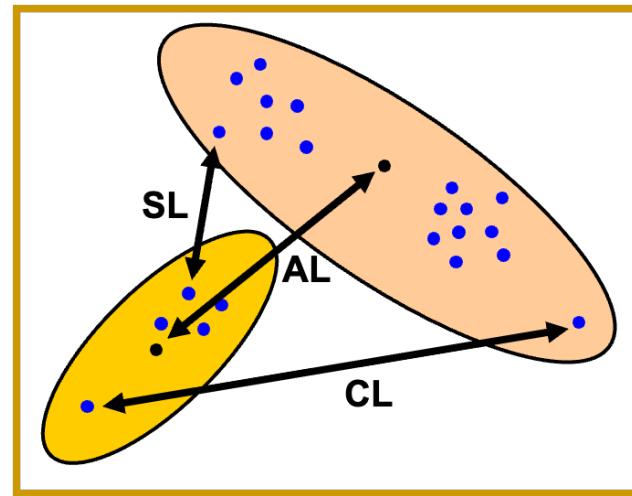
Similarity between clusters



- **Single linkage:** $d(X_i, X_j) = \min(d(a, b) : a \in X_i, b \in X_j)$
- **Complete linkage:** $d(X_i, X_j) = \max(d(a, b) : a \in X_i, b \in X_j)$
- **Average linkage:** $d(X_i, X_j) = \frac{1}{|X_i||X_j|} \sum_{a \in X_i} \sum_{b \in X_j} d(a, b)$

Hierarchical clustering

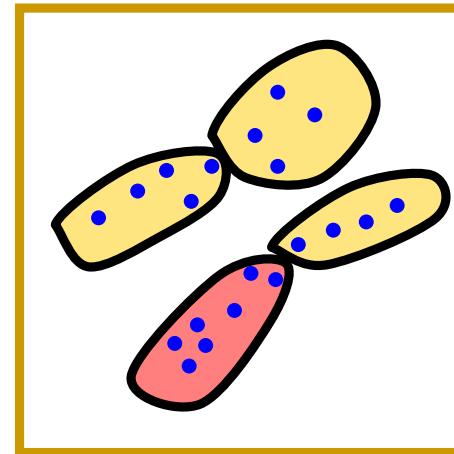
Similarity between clusters



- **Ward's method:** $d(X_i, X_j) = \|X_i - X_j\|^2$
(total variance)
- **Centroid distance:** $d(X_i, X_j) = \|c_s - c_t\|$

Hierarchical clustering

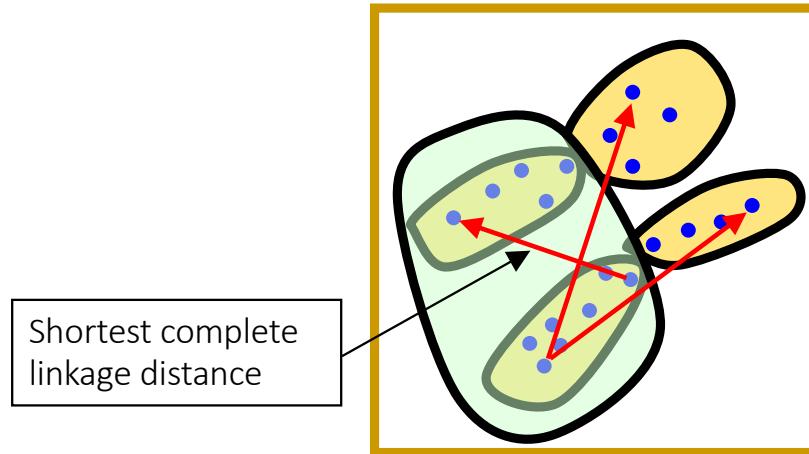
Similarity between clusters



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

Hierarchical clustering

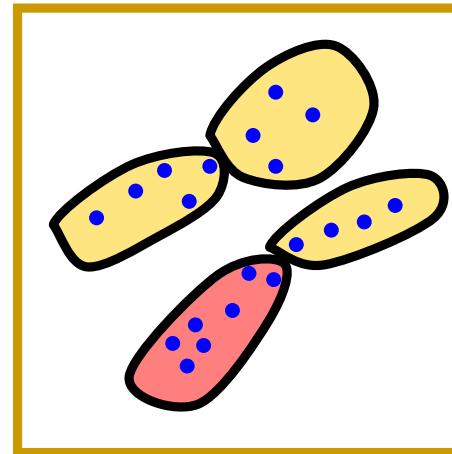
Similarity between clusters



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

Hierarchical clustering

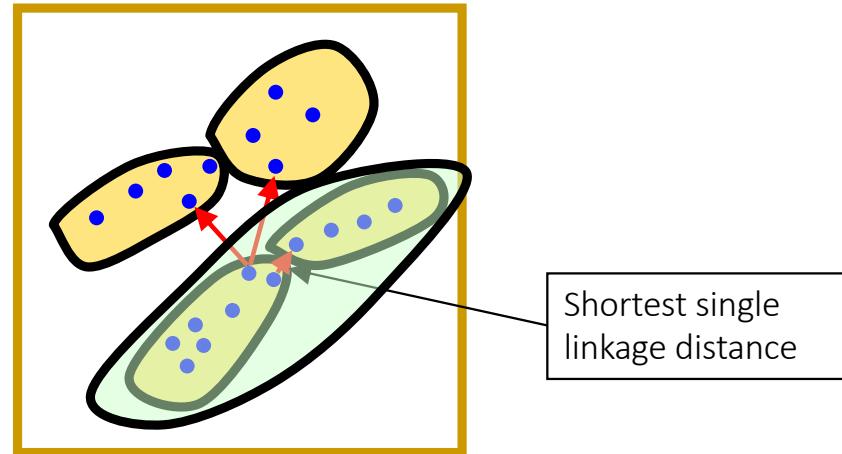
Similarity between clusters



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

Hierarchical clustering

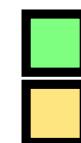
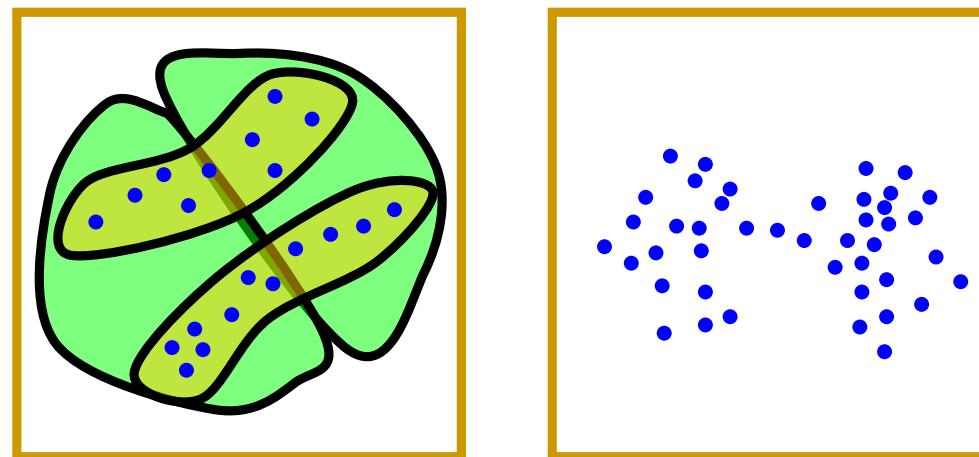
Similarity between clusters



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

Hierarchical clustering

Similarity between clusters

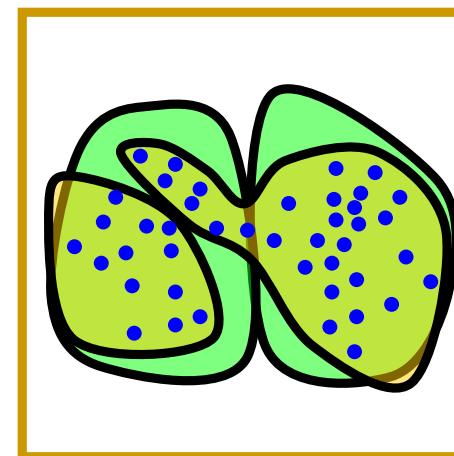
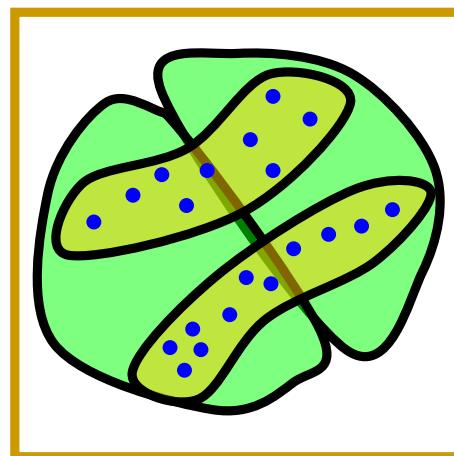


complete linkage
single linkage

Hierarchical clustering

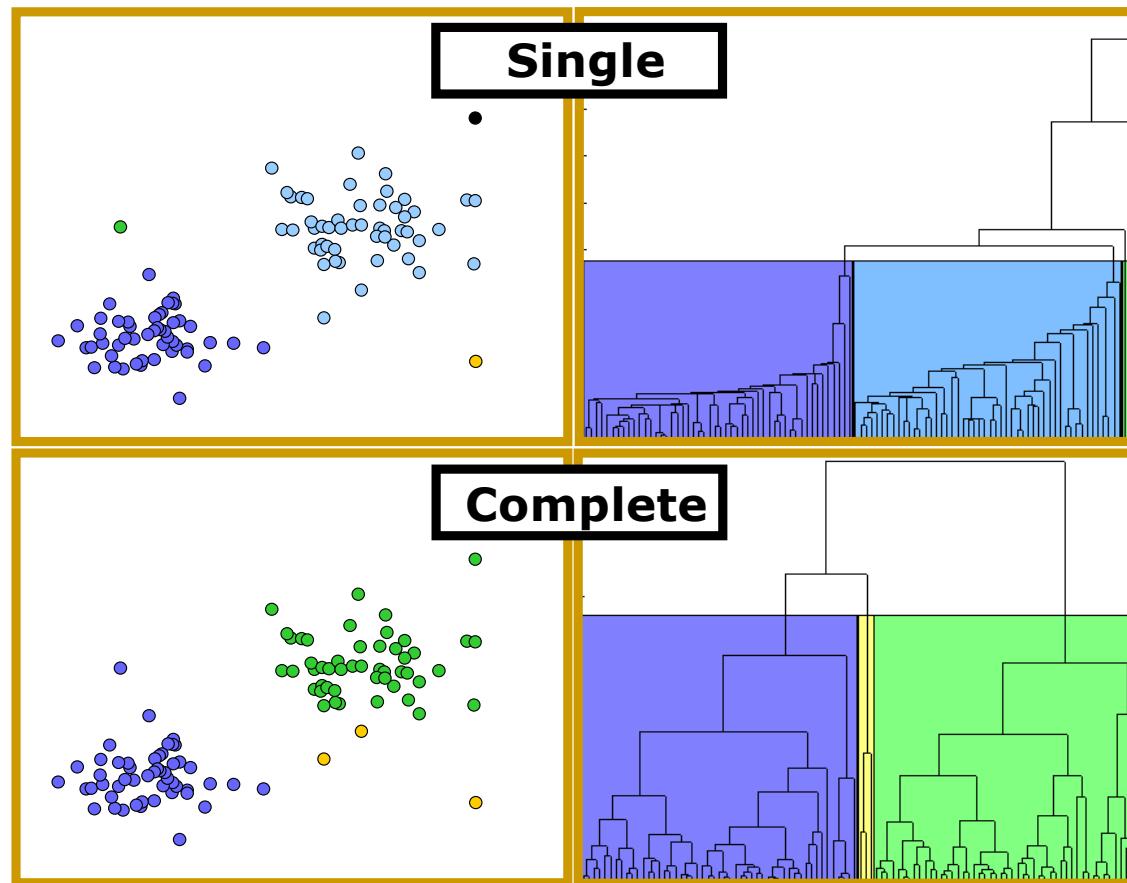
Similarity between clusters

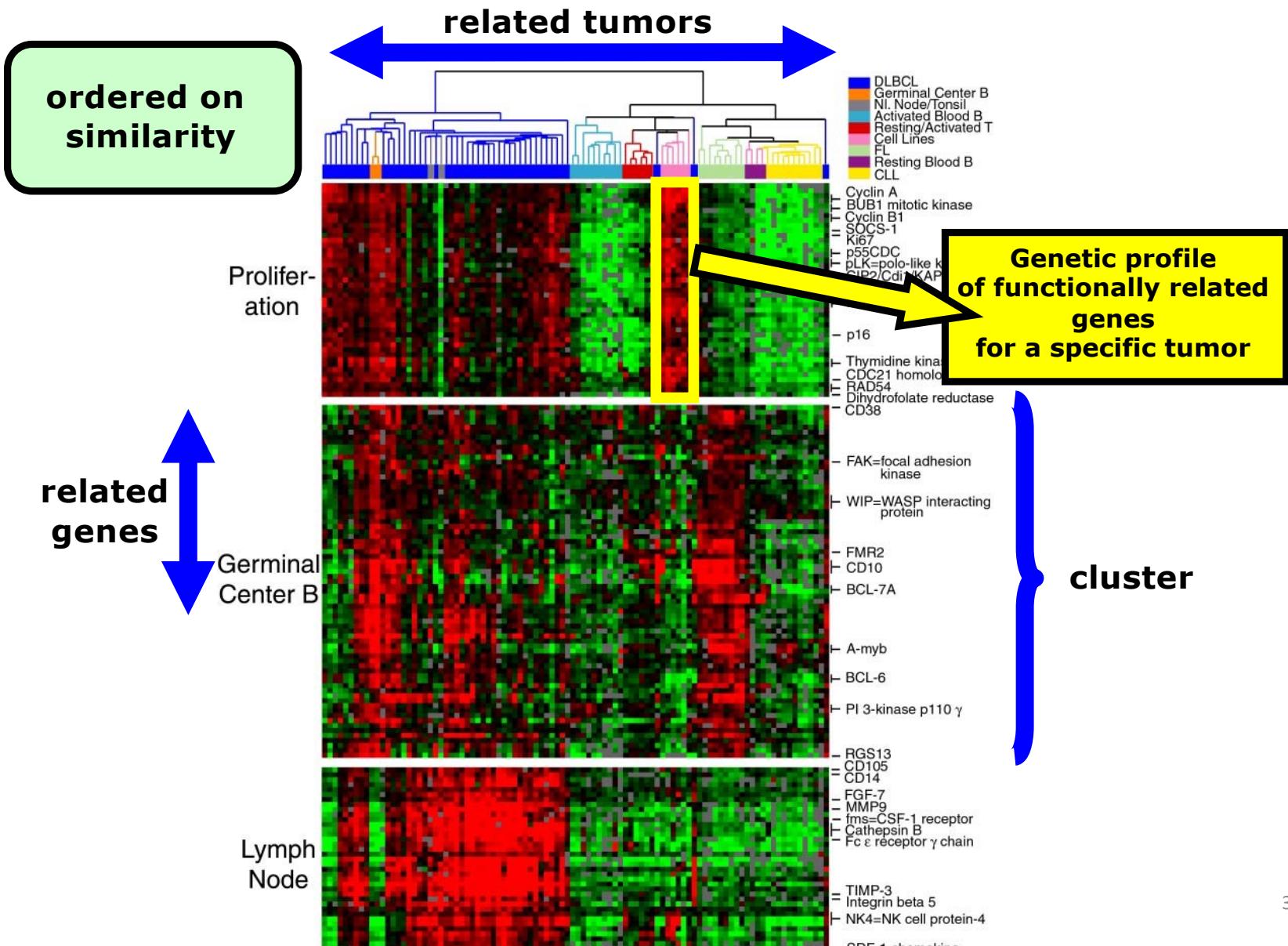
- Single linkage → long and “loose” clusters
- Complete linkage → compact clusters



complete linkage
single linkage

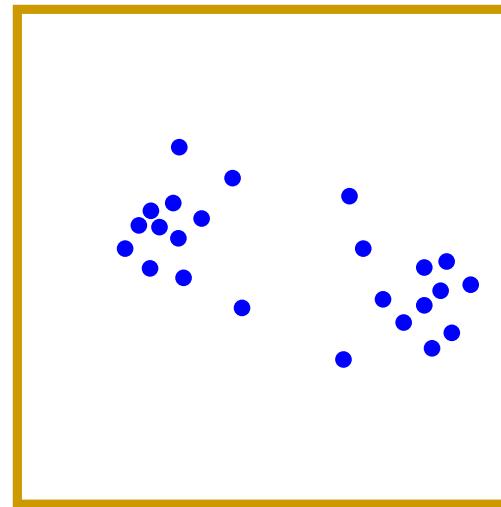
SL vs CL: Outlier influences



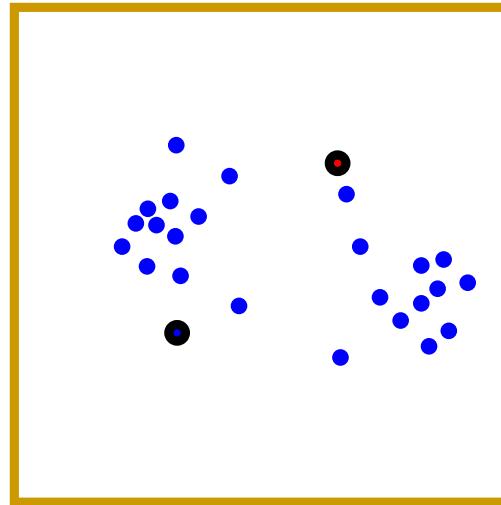


k-means clustering

k -Means clustering

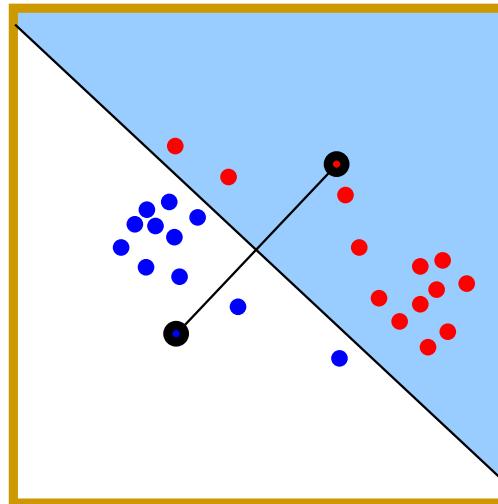


k -Means clustering



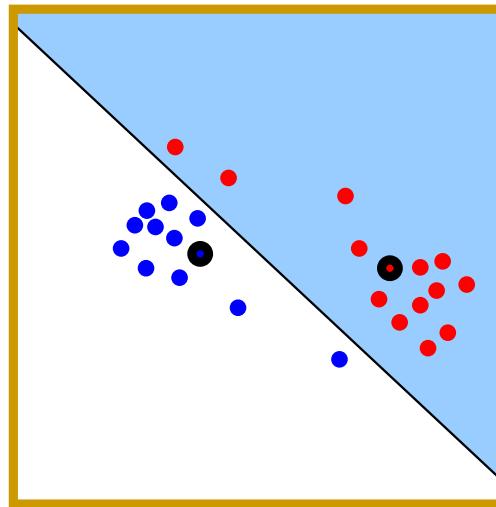
Choose randomly 2 prototypes

k -Means clustering



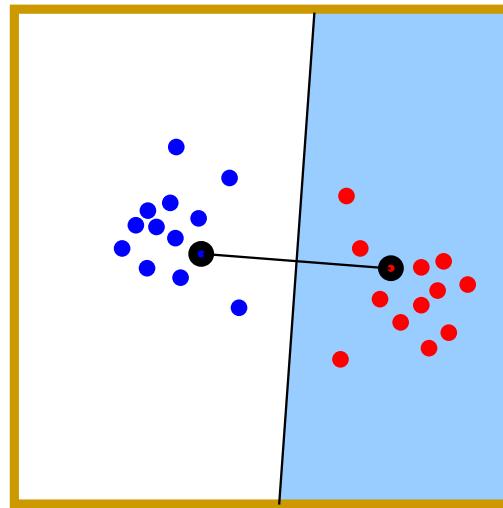
Assign objects to closest prototype
Blue area: cluster 1
White area: cluster 2

k -Means clustering



Calculate new cluster prototypes
By averaging objects

k -Means clustering

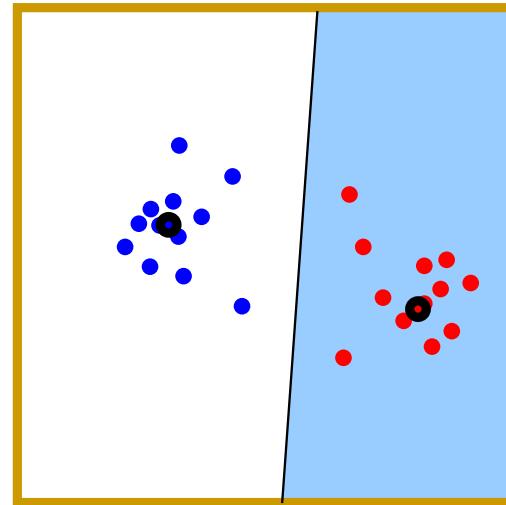


Re-assign objects to closest prototype

Blue area: cluster 1

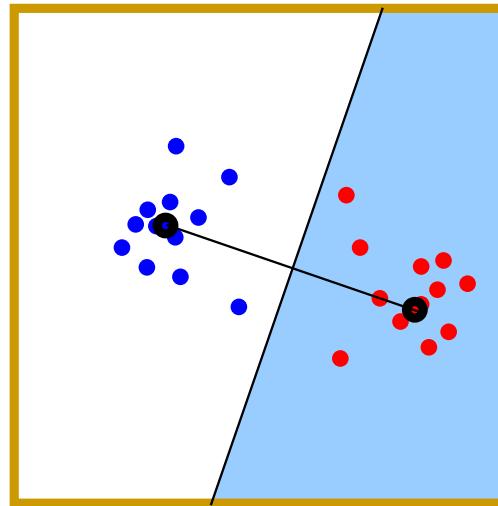
White area: cluster 2

k -Means clustering



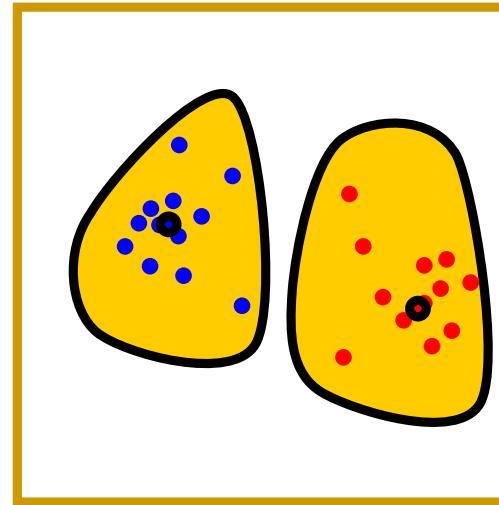
Re-calculate new cluster prototypes

k -Means clustering



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

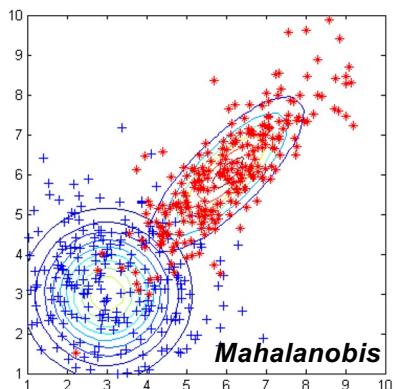
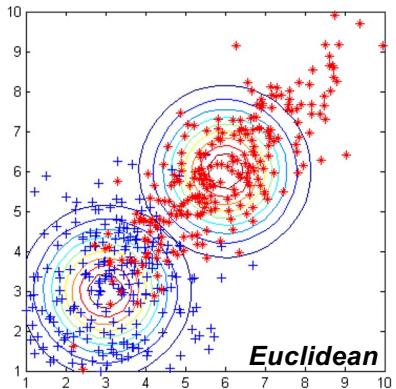
k -Means clustering



Establish clusters

k -Means clustering

Overview

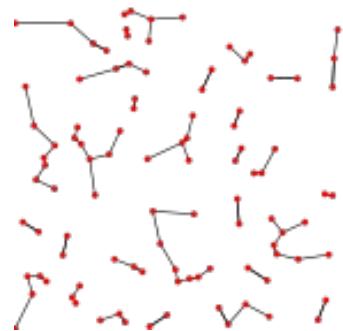


- k -Means
 - Fixed number of clusters (need to know a priori)
 - Choice of distance measure
 - Prototype choice
- Distance measure
 - Euclidean: Round clusters
 - Mahalanobis: Elongated clusters
- Prototype choice
 - Point
 - Line etc.
- Number of clusters
 - Validate clustering!

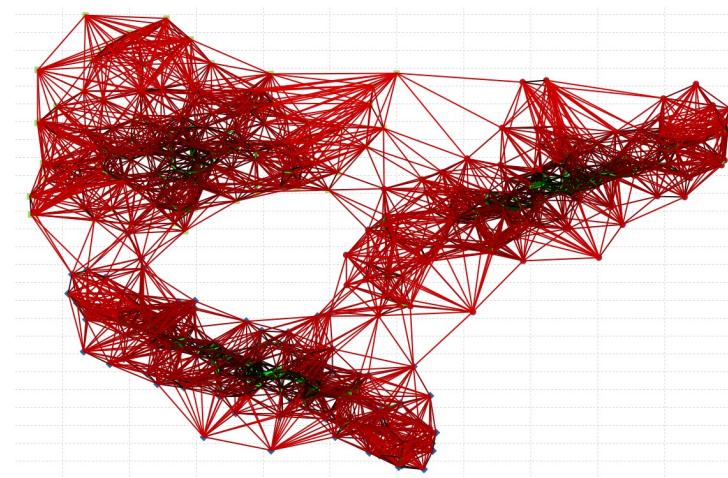
Graph-based Clustering

Graph-based clustering

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

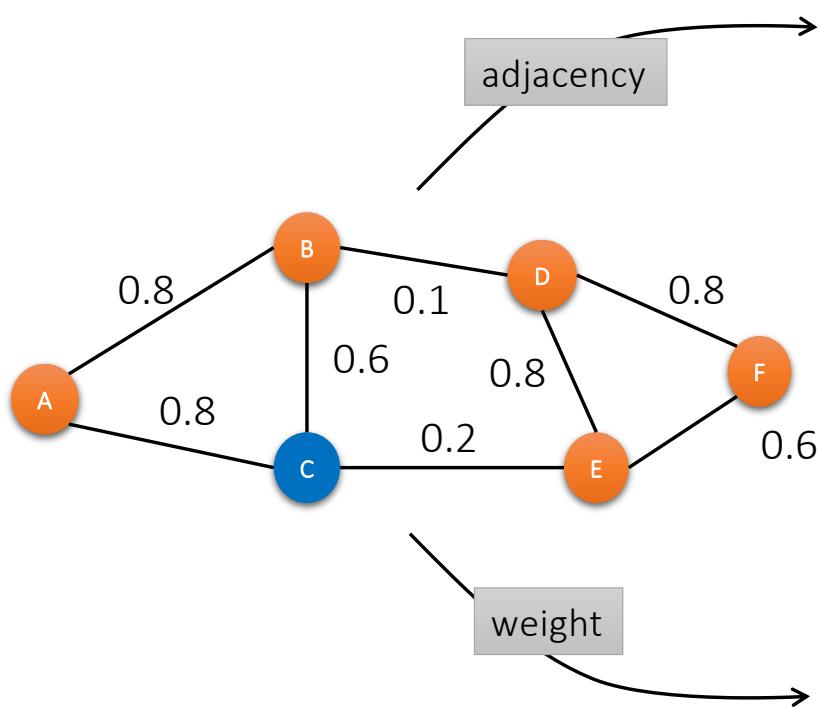


$k=1$



$k=20$

Graphs, adjacency and weight matrices

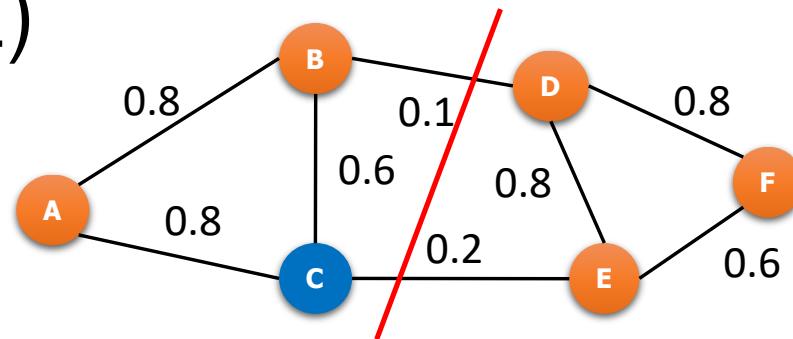


$$A = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \left(\begin{array}{cccccc} 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 1 & 0 \end{array} \right) \end{matrix}$$

$$W = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \left(\begin{array}{cccccc} 0 & 0.8 & 0.8 & 0 & 0 & 0 \\ 0.8 & 0 & 0.6 & 0.1 & 0 & 0 \\ 0.8 & 0.6 & 0 & 0 & 0.2 & 0 \\ 0 & 0.1 & 0 & 0 & 0.8 & 0.8 \\ 0 & 0 & 0.2 & 0.8 & 0 & 0.6 \\ 0 & 0 & 0 & 0.8 & 0.6 & 0 \end{array} \right) \end{matrix}$$

Spectral clustering (1)

- Minimise normalised cut
- Normalised cut between two clusters C_1 and C_2 :



$$NC(C_1, C_2) = \frac{\text{cut}(C_1, C_2)}{\text{assoc}(C_1, V)} + \frac{\text{cut}(C_2, C_1)}{\text{assoc}(C_2, V)} = 2 - \left(\frac{\text{assoc}(C_1, C_1)}{\text{assoc}(C_1, V)} + \frac{\text{assoc}(C_2, C_2)}{\text{assoc}(C_2, V)} \right)$$

- $\text{cut}(C_1, C_2)$ = weight of links between C_1 and C_2
- $\text{cut}(C_2, C_1)$ = same
- $\text{assoc}(C_1, V)$ = total weight of links from nodes in C_1 to entire graph
- $\text{assoc}(C_2, V)$ = total weight of links from nodes in C_2 to entire graph

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL. 22, NO. 8, AUGUST 2000

Normalized Cuts and Image Segmentation

Jianbo Shi and Jitendra Malik, Member, IEEE

Louvain method for community detection

Maximize modularity

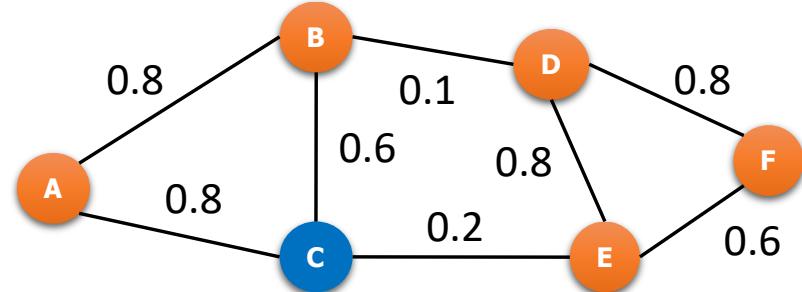
$$Q = \frac{1}{2m} \sum_{ij} \left[A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j)$$

Q_{ij}

Observed edges in cluster c

Expected edges in cluster c

i, j	nodes
c_i, c_j	community of node i and j
A_{ij}	weight edge node i and j
k_i, k_j	total weight of edges of nodes i and j
$\delta(c_i, c_j)$	1 if $c_i = c_j$, zero otherwise
m	total weight of edges

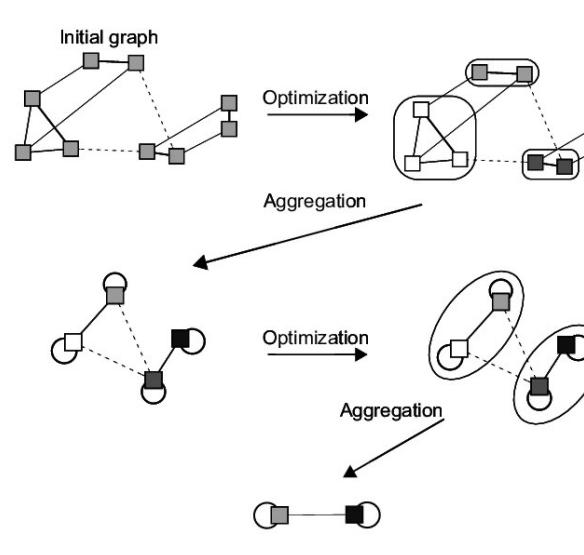


Louvain method for community detection

Modularity optimized in two phases, applied iteratively

Phase 1: Optimization (communities); moving nodes from own community to community of neighboring node (*improve Q*)

Phase 2: Aggregation; building new network by aggregating nodes in detected communities



Outline

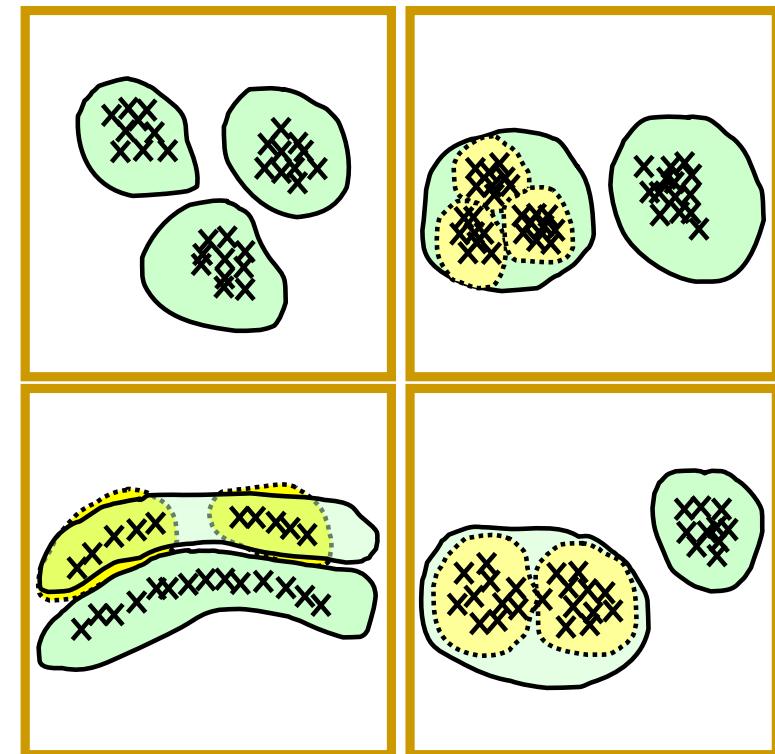
- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering

• Cluster validation

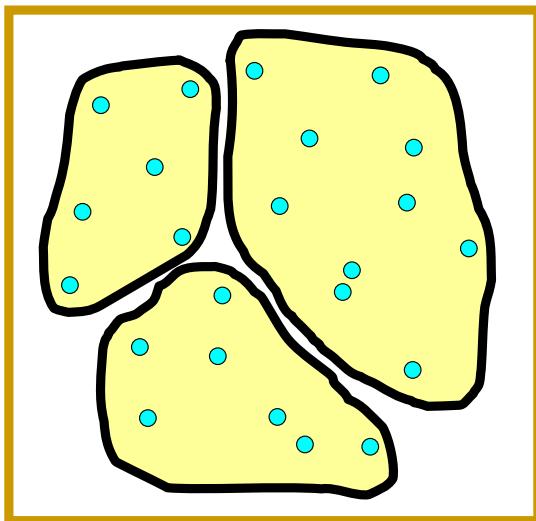
- scRNA-seq clustering
- Annotating clusters

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 Validation



- Cluster tendency

Clustering **IMPOSES** structure even though data may not possess it

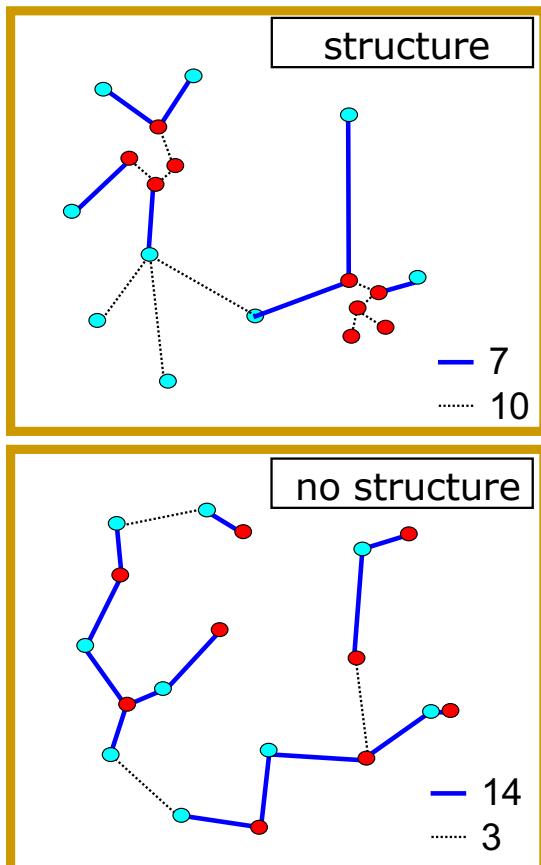
Aim: Test whether data possesses structure

- Cluster validity

Choices impose restrictions on for example shape

Aim: Quantitative evaluation of the clustering results

Test for spatial randomness



- Test

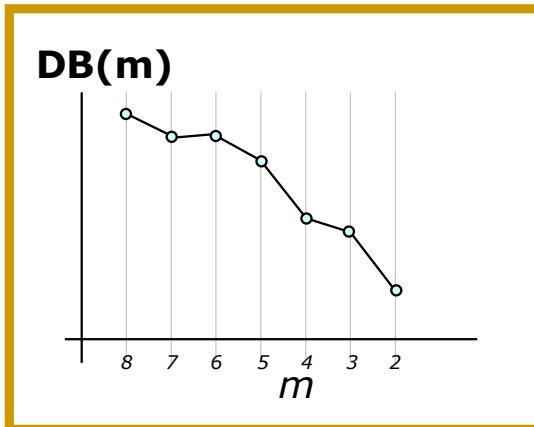
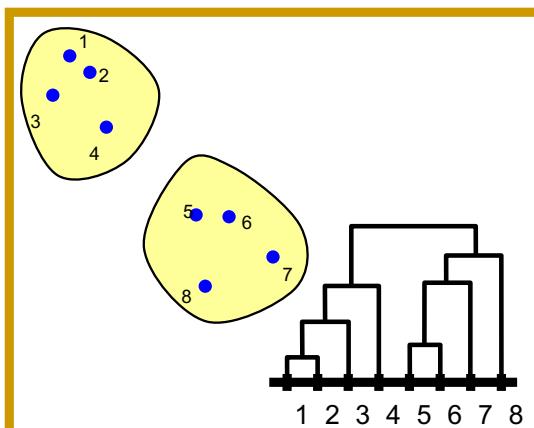
If data (●) clusters frequently with random data (○)
then data structureless

- Approach

- Generate random vectors (\mathbf{Y}) uniformly over observed region of data (\mathbf{X})
- Find MST (single linkage HC) of $\mathbf{X} \vee \mathbf{Y}$
- Determine number of edges q that connect vectors of \mathbf{X} with \mathbf{Y}
- If \mathbf{X} contains clusters q should be small!

(multiple random vs random measurements gives likelihood for q)

Davis-Bouldin index



- Test

Select specific clustering according to a criteria

For example: Davis-Bouldin index

- DB index

For a specific clustering m , $\mathbf{DB}(m)$: Average similarity of a cluster with its most similar cluster

- Approach

Goal: Clusters to have minimal similarity

Seek: Clustering that minimize $DB(m)$ wrt m

Davis-Bouldin index

- Similarity cluster C_i and C_j

$$D_{i,j} = \frac{(\bar{d}_i + \bar{d}_j)}{\|\mu_i - \mu_j\|}$$

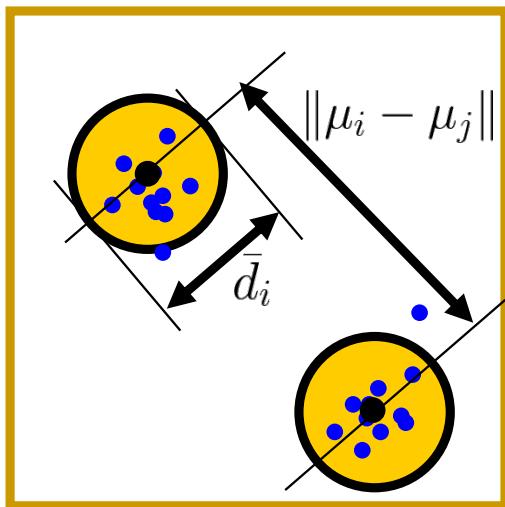
• \bar{d}_i : average distance within cluster i , μ_i : centroid of cluster i

- Most similar cluster to C_i

$$R_{i,j} = \max_{j \neq i} \{D_{i,j}\}$$

- DB index

$$DB = \frac{1}{k} \sum_{k=1}^k R_{i,j}$$



Silhouette score

- Measure similarity of object to its own cluster

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

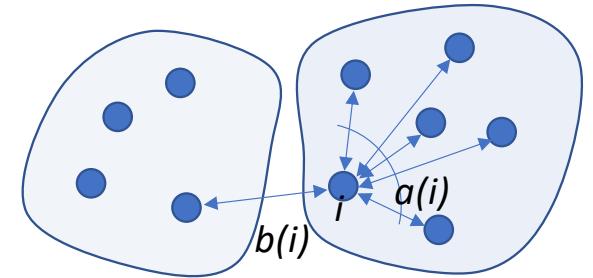
with $a(i)$ being average distance to all objects in same cluster and $b(i)$ being closest object from all other clusters:

$$a(i) = \frac{1}{|C_i|} \sum_{\forall j} d(x_i, x_j) \quad b(i) = \min_{\forall j, j \notin C_i} d(x_i, x_j)$$

$-1 \leq s(i) \ll 1$; $s(i)$ is close to 1, if $a(i) \ll b(i)$; average distance within cluster much smaller than nearest objects

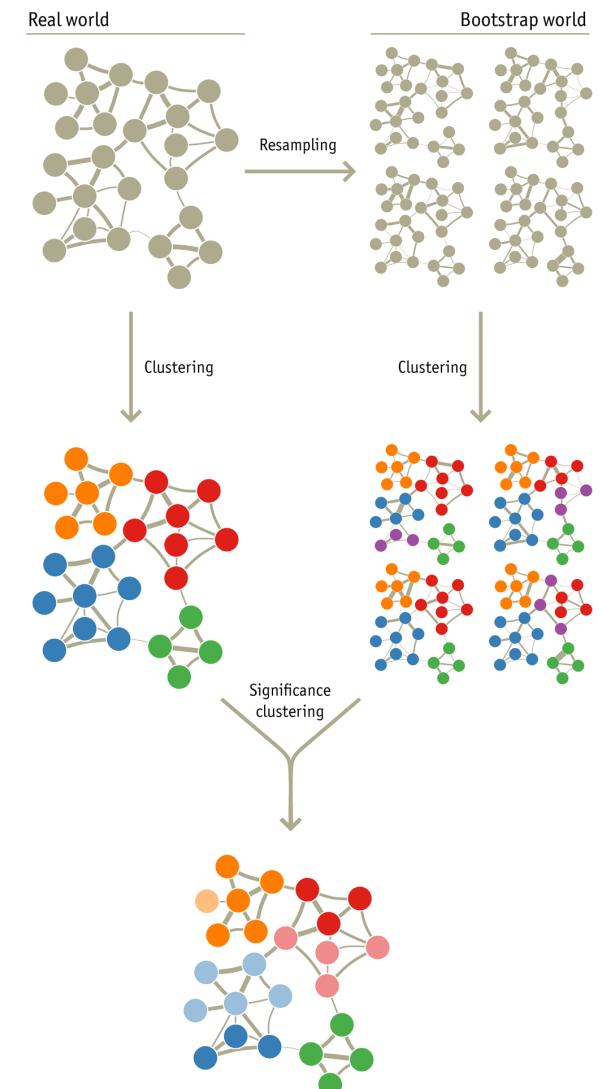
- Silhouette score is average of all these similarities

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



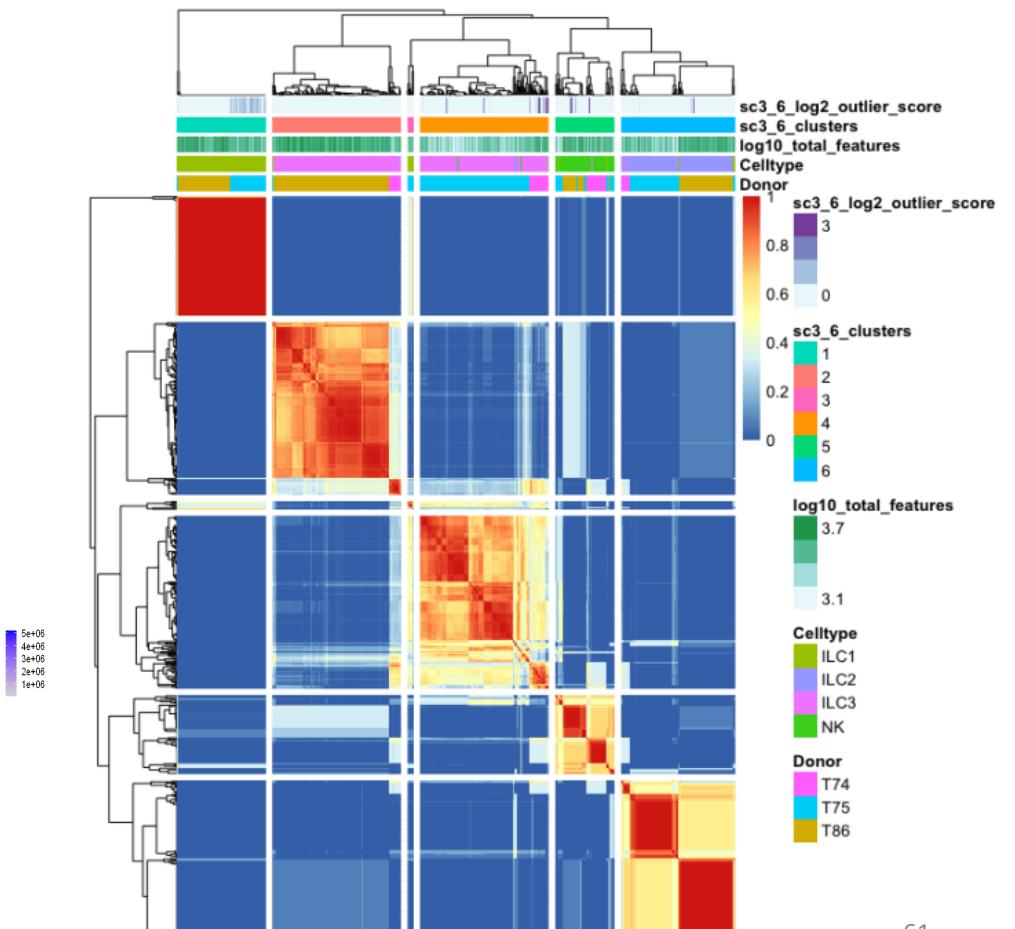
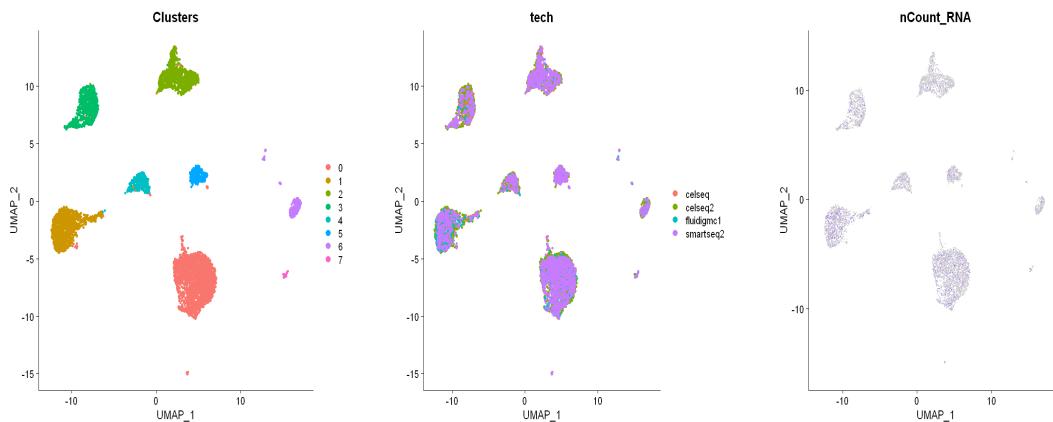
Bootstrapping

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



Always check QC data

- Is what your splitting mainly related to batches, qc-measures (especially detected genes)?



Outline

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

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

Comparing different cluster labels

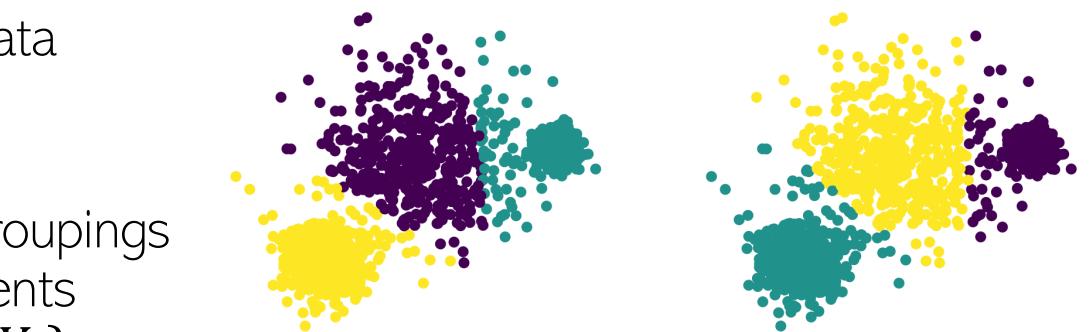
Adjusted Rand Index (ARI)

Measure of the similarity between two data clusterings

Given a set S of n elements, and two groupings or partitions (clusterings) 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	

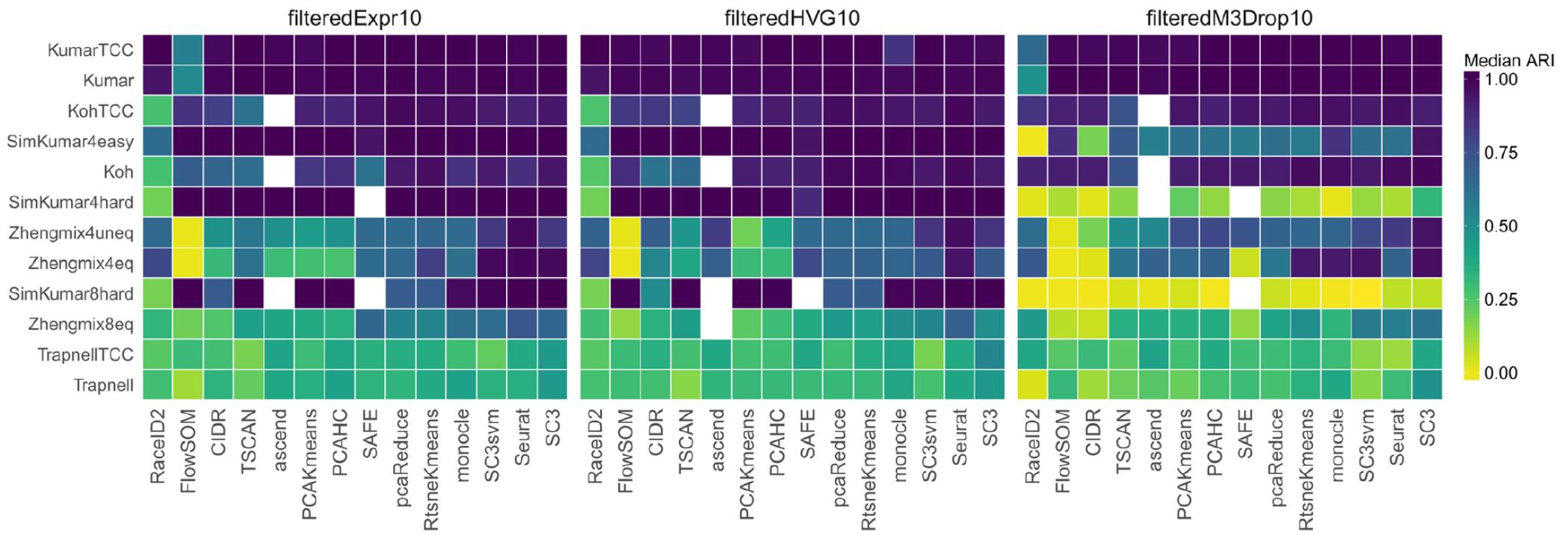
$$n_{ij} = |X_i \cap Y_j|$$



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

$-1 \leq \text{ARI} \leq 1$ (1: perfect overlap, 0: random overlap)

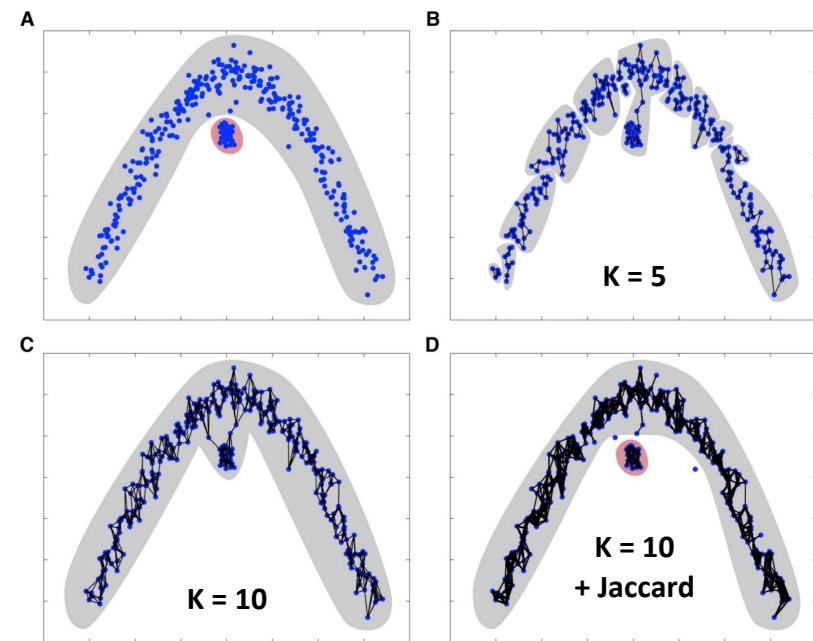
Benchmarking scRNA-seq clustering methods



Row: different data set, Panel: different gene filtering, Column: different clustering method

Standard clustering approach

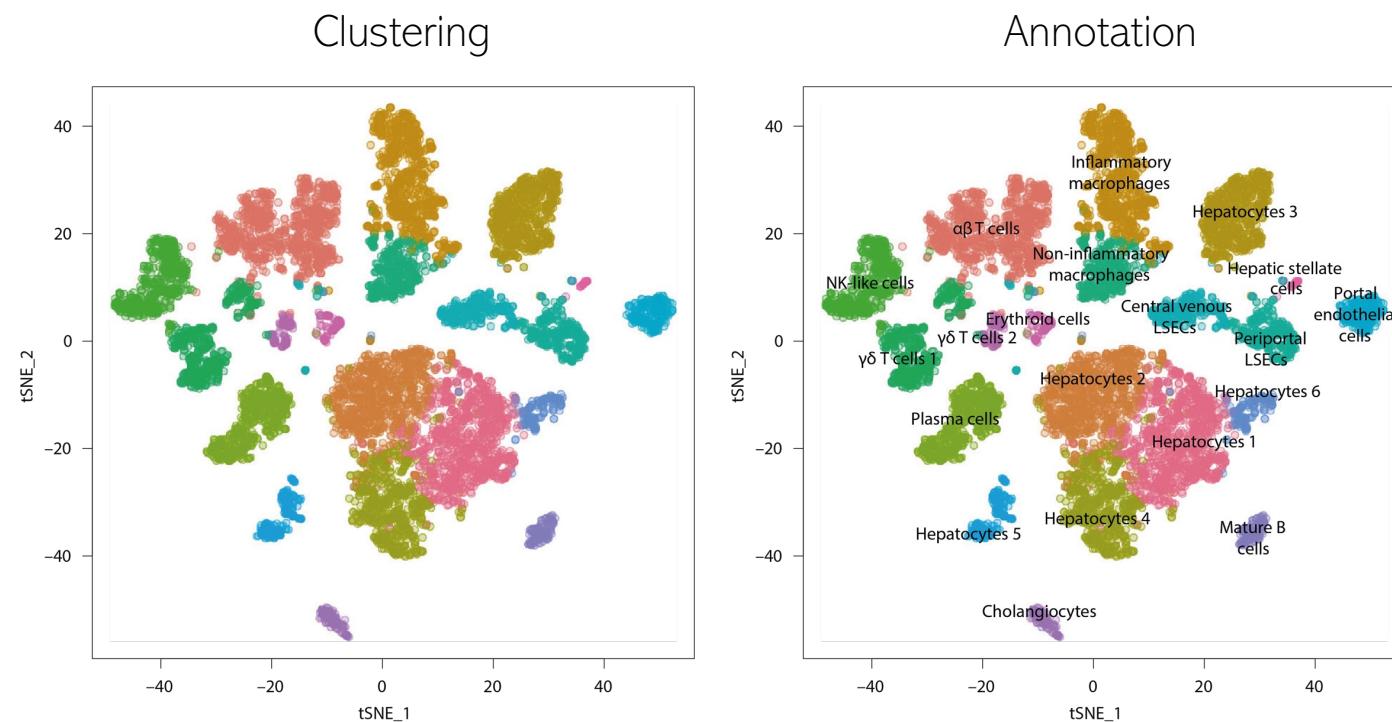
1. Select highly variable genes (~1000-5000 genes)
2. Reduce dimensions using PCA (~30-50 dimensions)
3. Construct kNN graph (~15-20 neighbors)
4. Louvain/Leiden community detection



Outline

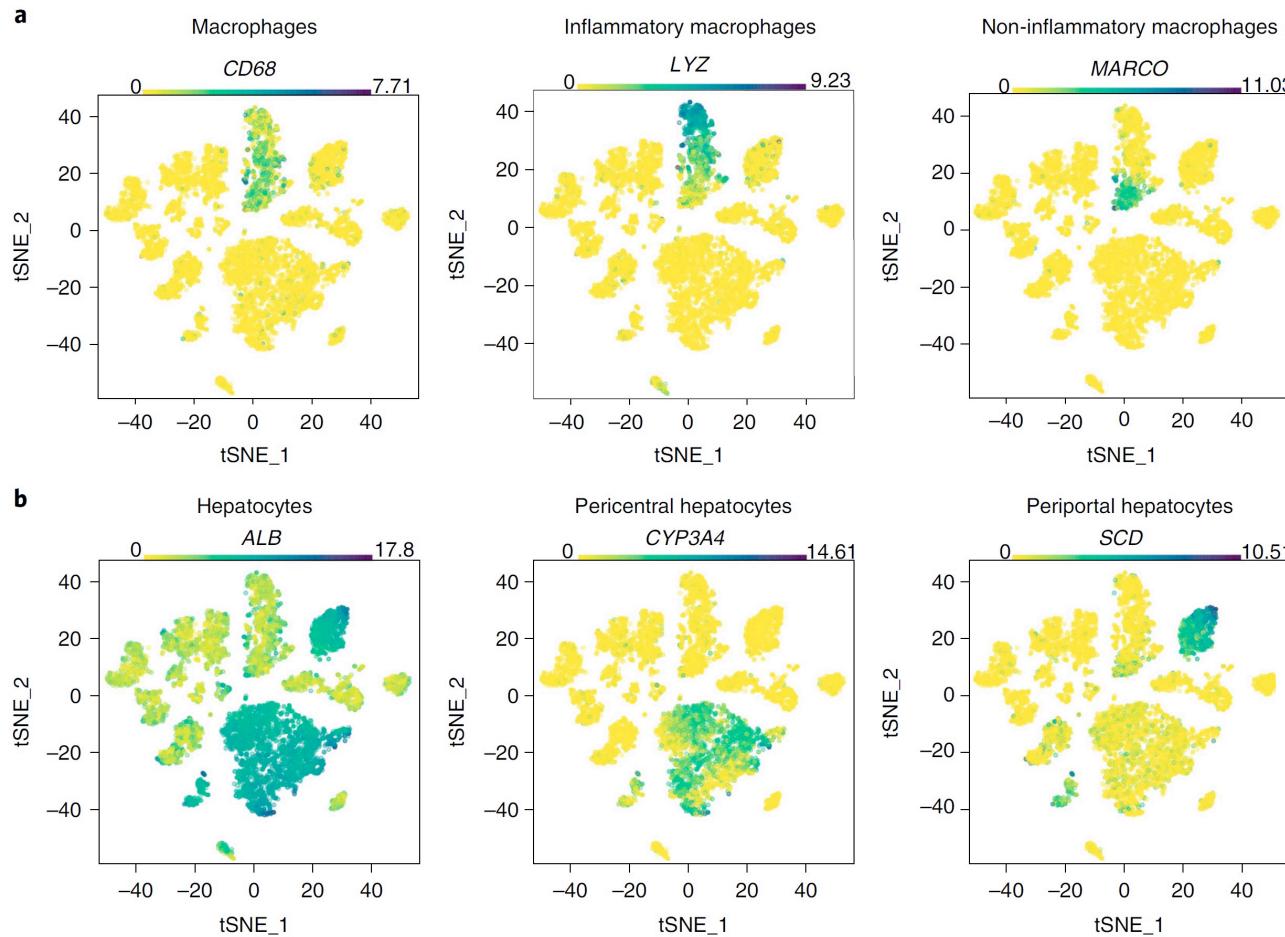
- Introduction to clustering
 - Hierarchical clustering
 - k -Means clustering
 - Graph-based clustering
- Cluster validation
- scRNA-seq clustering
- **Annotating clusters**

From clusters to annotations



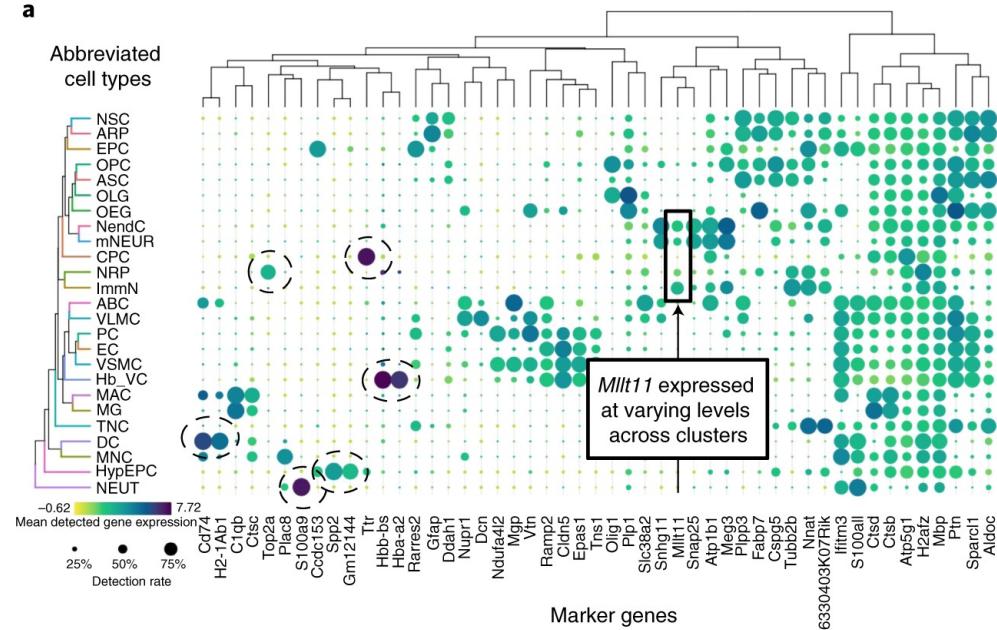
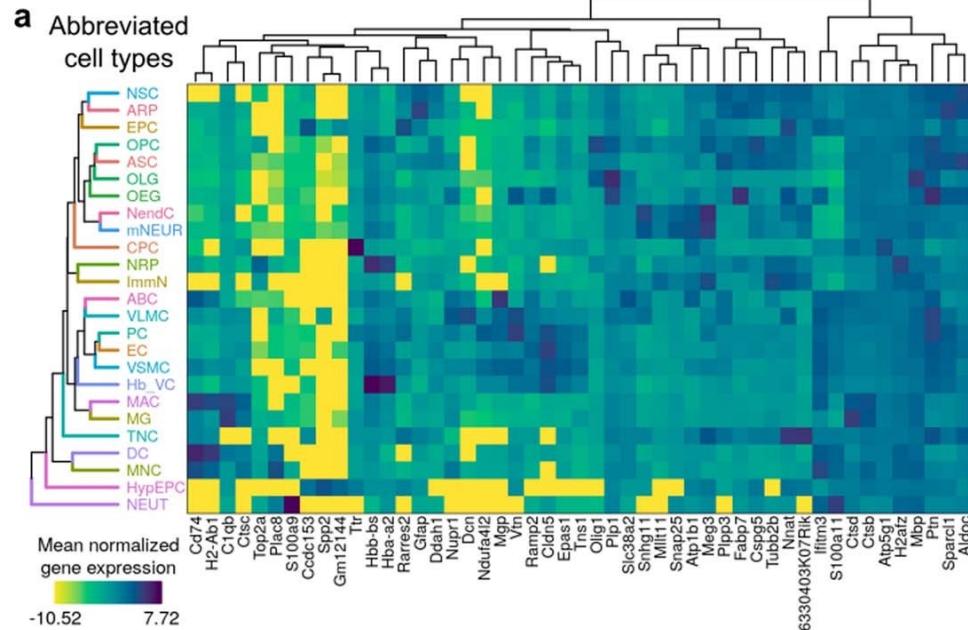
Gene expression overlay

Easy



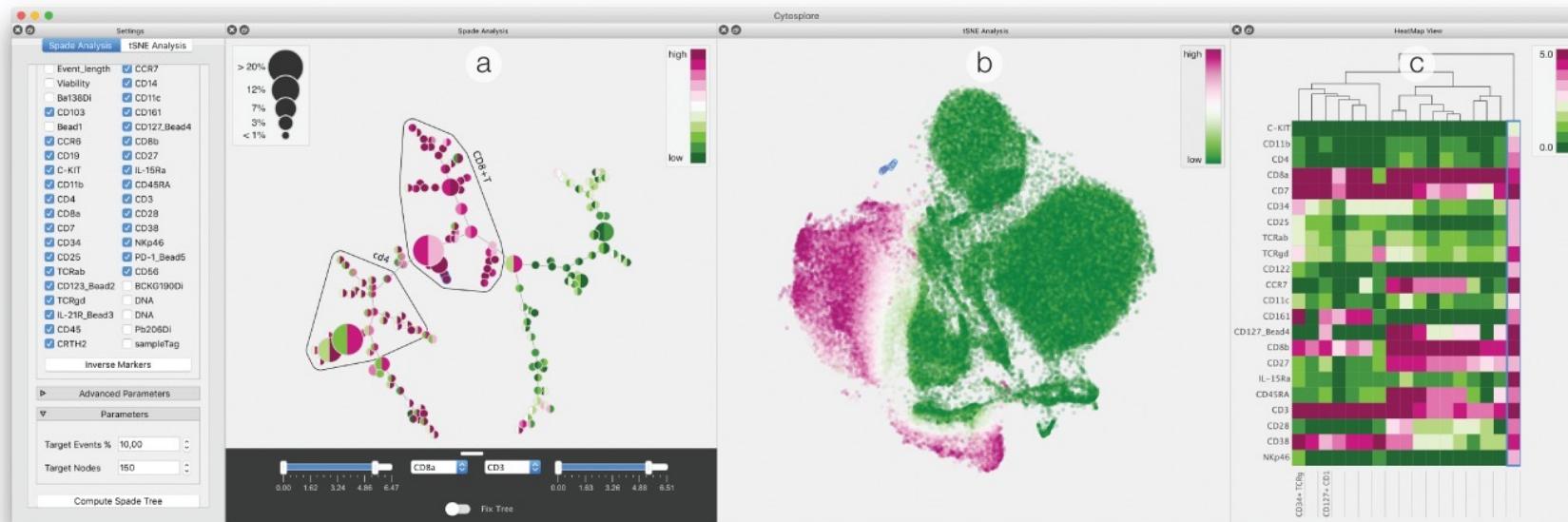
Challenging

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

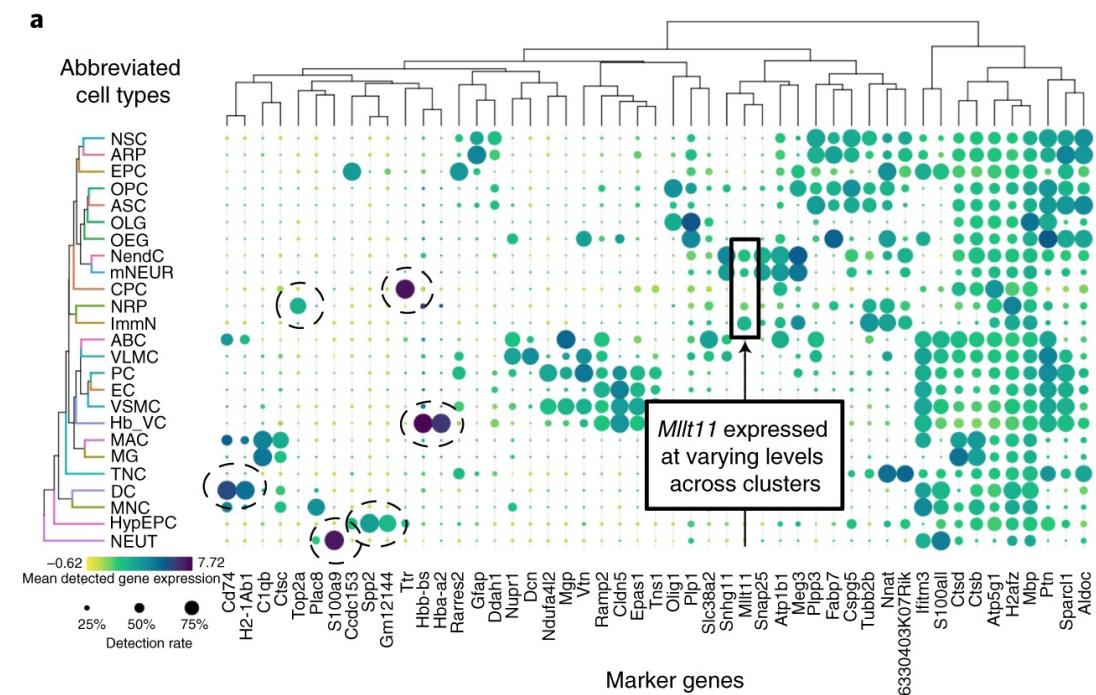
- “Expert knowledge”
- Literature
- Other scRNA-seq data
- **Marker databases: PangaloDB, CellMarker,...**

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 markers to find functional information that may help identify the cell type



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?

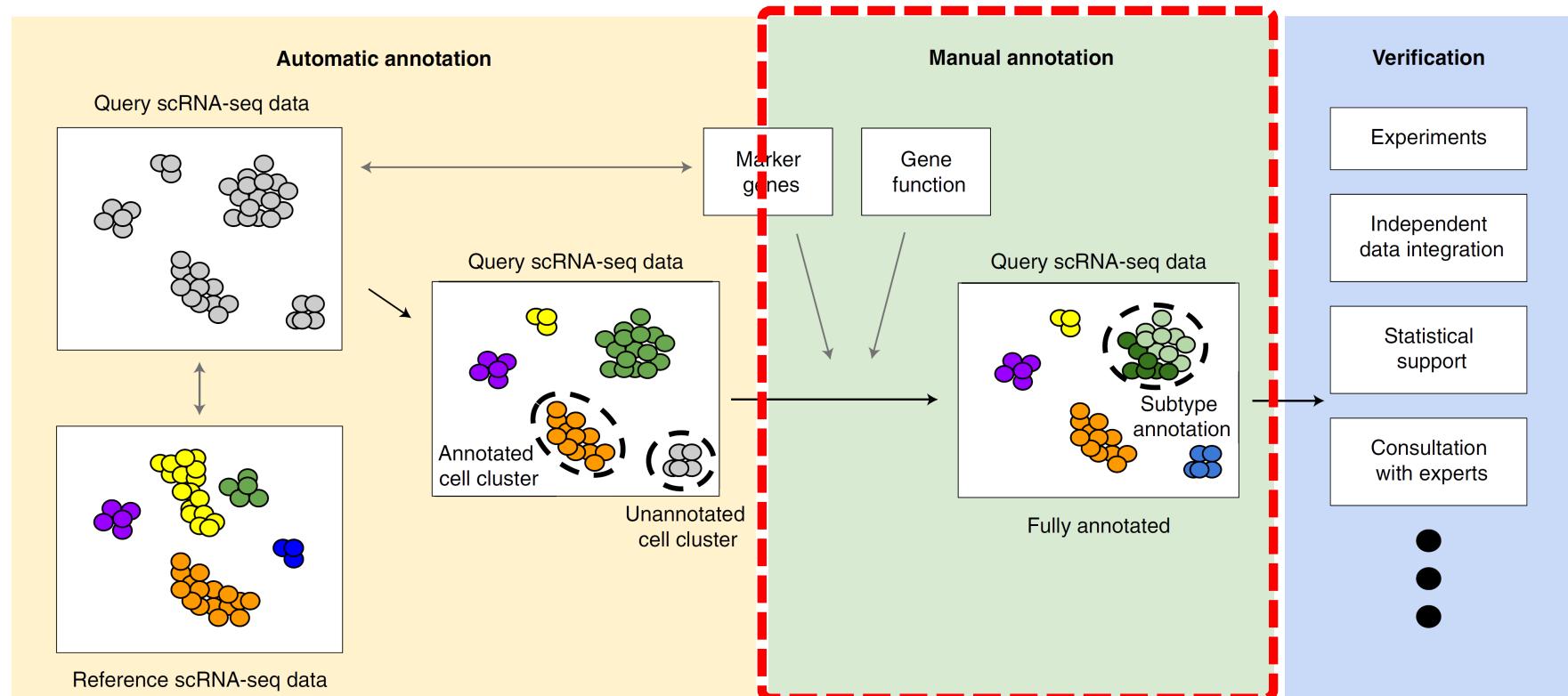
Miller et al. (eLife 2020)

<https://portal.brain-map.org/explore/classes/nomenclature>

Cell set nomenclature	Hierarchical organization of cell types defined for this taxonomy	Cell type alias	Cell type label	Cell type accession ID
A: Alias: GABAergic Label: Neuron 1-23 ID: CS1601040051		Smad3 12 Ndnf Car4 24 Ndnf Cxd14 30 Igtp 10 Vip Gpc3 45 Vip Chat 46 Vip Parml 38 Vip Mybpc1 24 Vip Snrg 13 Snog 9 Sst Chodl 41 Sst Cdk6 14 Sst Cbln4 66 Sst Myh8 38 Sst Tacstd2 12 Sst Th 16 Pvalb Cpne5 14 Pvalb Tpbpg 12 Pvalb Obox3 16 Pvalb Gpc3 54 Pvalb Rspo2 21 Pvalb Wt1 46 Pvalb Tacr3 63 L2 Ngb 16 L2/3 Ptgs2 92 L4 Arf5 31 L4 Scnn1a 61 L4 Cbx3 55 L5a Hes11b1 42 L5a Tcoag11 20 L5a Pde1c 12 L5 Ucma 12 L5a Batf3 57 L6a Car12 14 L6a Syt17 12 L5b Tph2 25 L5b Cd13 30 L6a Mgp 37 L6a Sia 53 L5b Chrat6 8 L6b Serpinb11 16 L6b Rgs12 13 Oligo Opalin 30 Oligo 96'Rik 7 OPC Pdgfra 22 Astro Aqp4 43 SMC My9 13 Endo Xdh 14 Micro Ctss 22	Neuron 1 Neuron 2 Neuron 3 Neuron 4 Neuron 5 Neuron 6 Neuron 7 Neuron 8 Neuron 9 Neuron 10 Neuron 11 Neuron 12 Neuron 13 Neuron 14 Neuron 15 Neuron 16 Neuron 17 Neuron 18 Neuron 19 Neuron 20 Neuron 21 Neuron 22 Neuron 23 Neuron 24 Neuron 25 Neuron 26 Neuron 27 Neuron 28 Neuron 29 Neuron 30 Neuron 31 Neuron 32 Neuron 33 Neuron 34 Neuron 35 Neuron 36 Neuron 37 Neuron 38 Neuron 39 Neuron 40 Neuron 41 Neuron 42 Non-neuron 1 Non-neuron 2 Non-neuron 3 Non-neuron 4 Non-neuron 5 Non-neuron 6 Non-neuron 7	CS1601040001 CS1601040002 CS1601040003 CS1601040004 CS1601040005 CS1601040006 CS1601040007 CS1601040008 CS1601040009 CS1601040010 CS1601040011 CS1601040012 CS1601040013 CS1601040014 CS1601040015 CS1601040016 CS1601040017 CS1601040018 CS1601040019 CS1601040020 CS1601040021 CS1601040022 CS1601040023 CS1601040024 CS1601040025 CS1601040026 CS1601040027 CS1601040028 CS1601040029 CS1601040030 CS1601040031 CS1601040032 CS1601040033 CS1601040034 CS1601040035 CS1601040036 CS1601040037 CS1601040038 CS1601040039 CS1601040040 CS1601040041 CS1601040042 CS1601040043 CS1601040044 CS1601040045 CS1601040046 CS1601040047 CS1601040048 CS1601040049
B: Alias: Pvalb Label: Neuron 17-23 ID: CS1601040054				
C: Alias: [blank] Label: Neuron 20-23 ID: CS1601040056				
D: Alias: L2/3 Label: Neuron 24-25 ID: CS1601040066				
E: Alias: Glutamatergic Label: Neuron 24-42 ID: CS1601040060				
F: Alias: [blank] Label: Neuron 24-42 ID: CS1601040066				
G: Alias: Non-neuronal Label: Non-neuron 1-7 ID: CS1601040070				



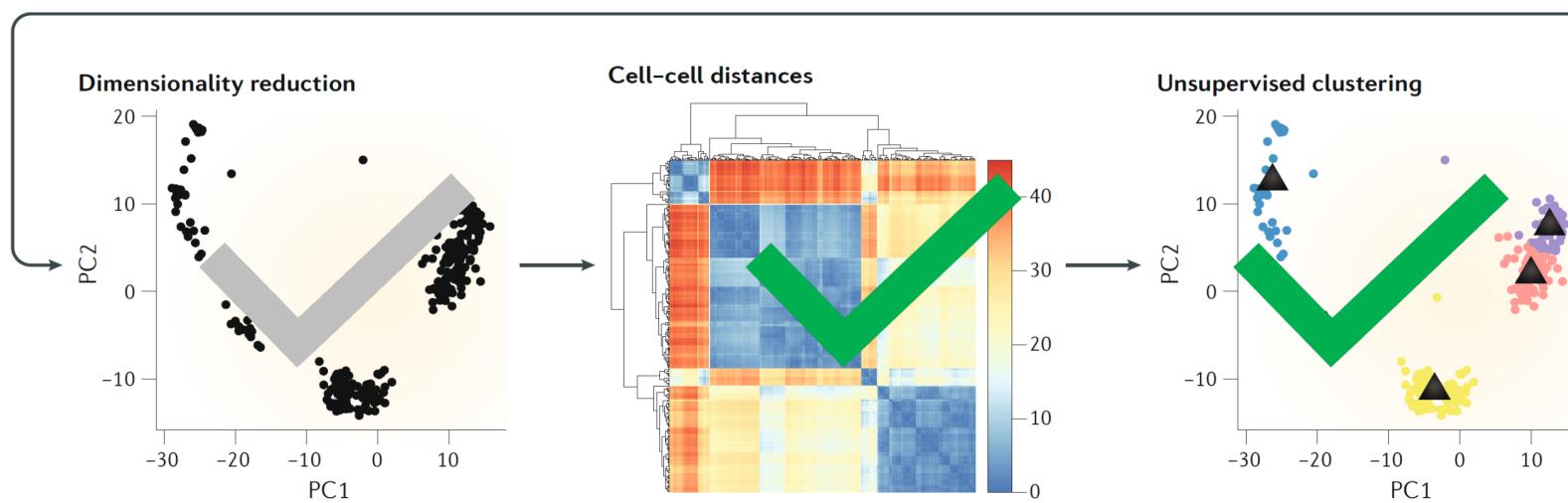
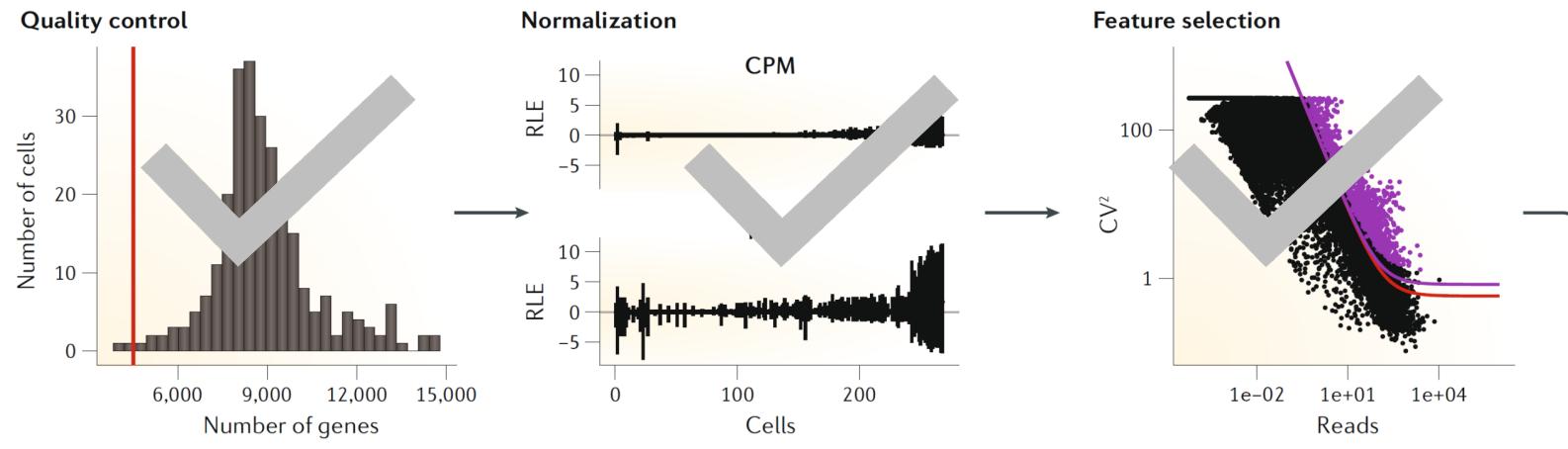
Tutorial: guidelines for annotating single-cell transcriptomic maps using automated and manual methods



Summary

- Start by identifying major well-known cell types (clearly defined, discrete cell clusters)
- Sometimes, it is useful to 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 relies on marker genes and expert knowledge

Summary



Challenges

- **Subjectivity:** what is a cell type
 - *Different parameters will give different results*
 - *Validation is important*
- **Scalability:** in the last few years the number of cells in scRNA-seq experiments has grown by several orders of magnitude from $\sim 10^2$ to $\sim 10^6$
 - *Computational efficiency*
 - *Visual exploration, crowding problem*

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>