

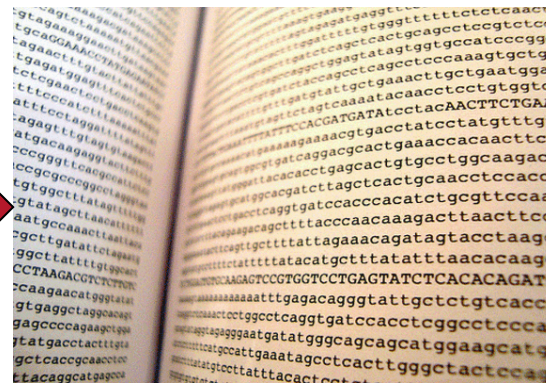
Clustering I

With application to gene-expression
profiling technology

Arjun Krishnan

Thanks to Kevin Wayne, Matt Hibbs, & SMD for a few of the slides

Why is expression important?

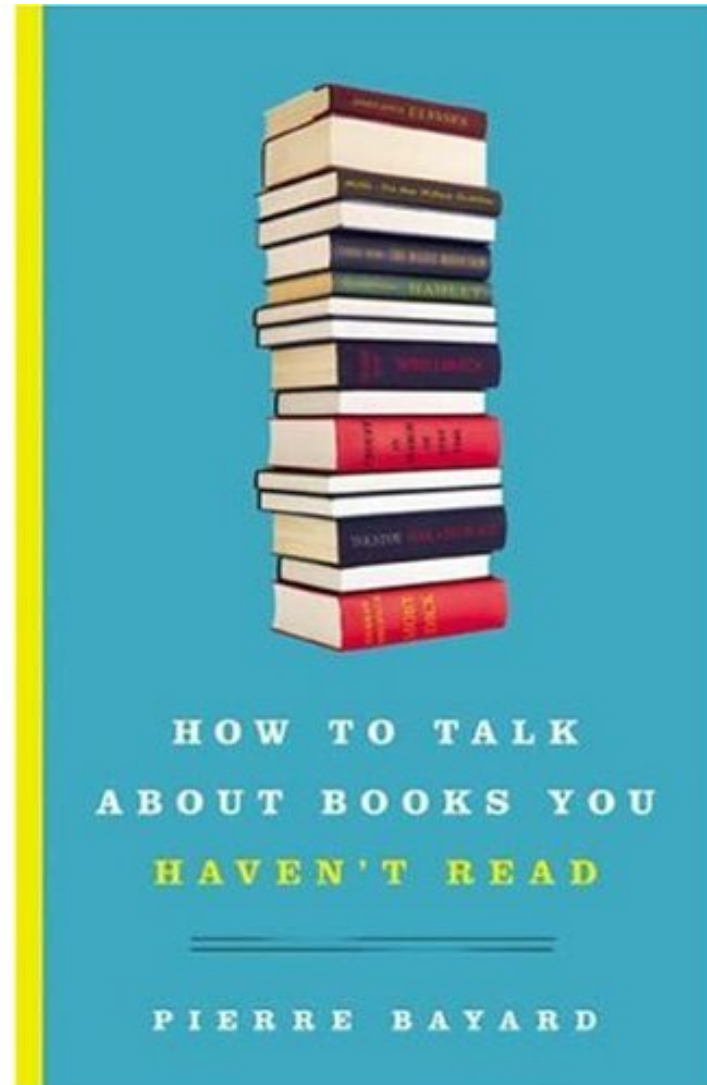


Understanding cellular and human biology

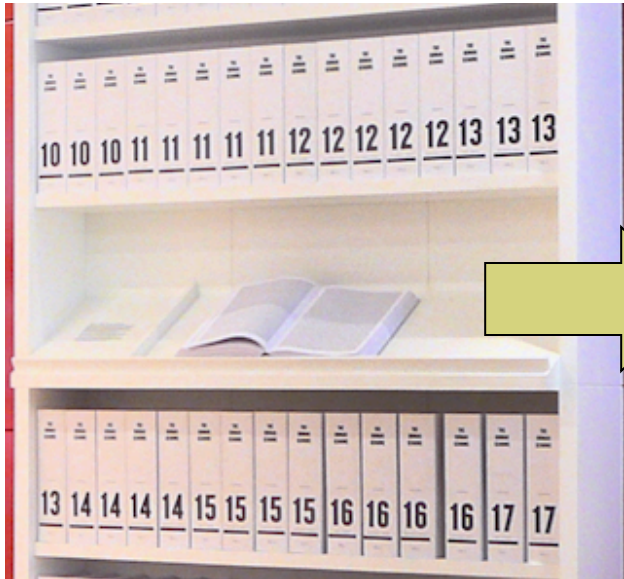


Understanding culture and social dynamics

Why is expression important?



Why is expression important?



Measure the activity of genes in various cellular conditions

Understanding cellular and human biology



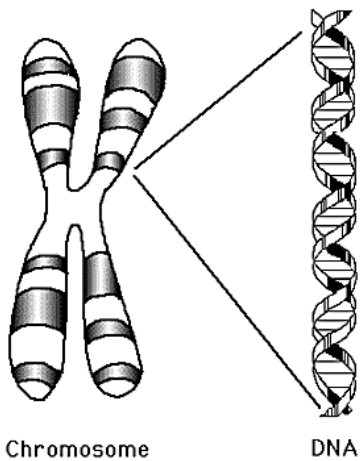
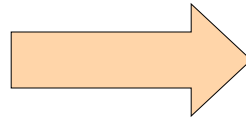
Measure the activity of people in various social instances

Understanding culture and social dynamics

Why is expression important?

Proteins

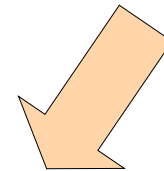
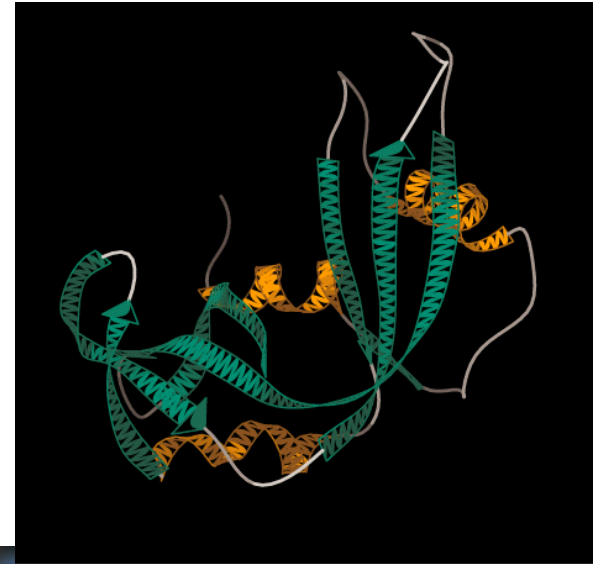
Gene
Expression



Chromosome

DNA

DNA



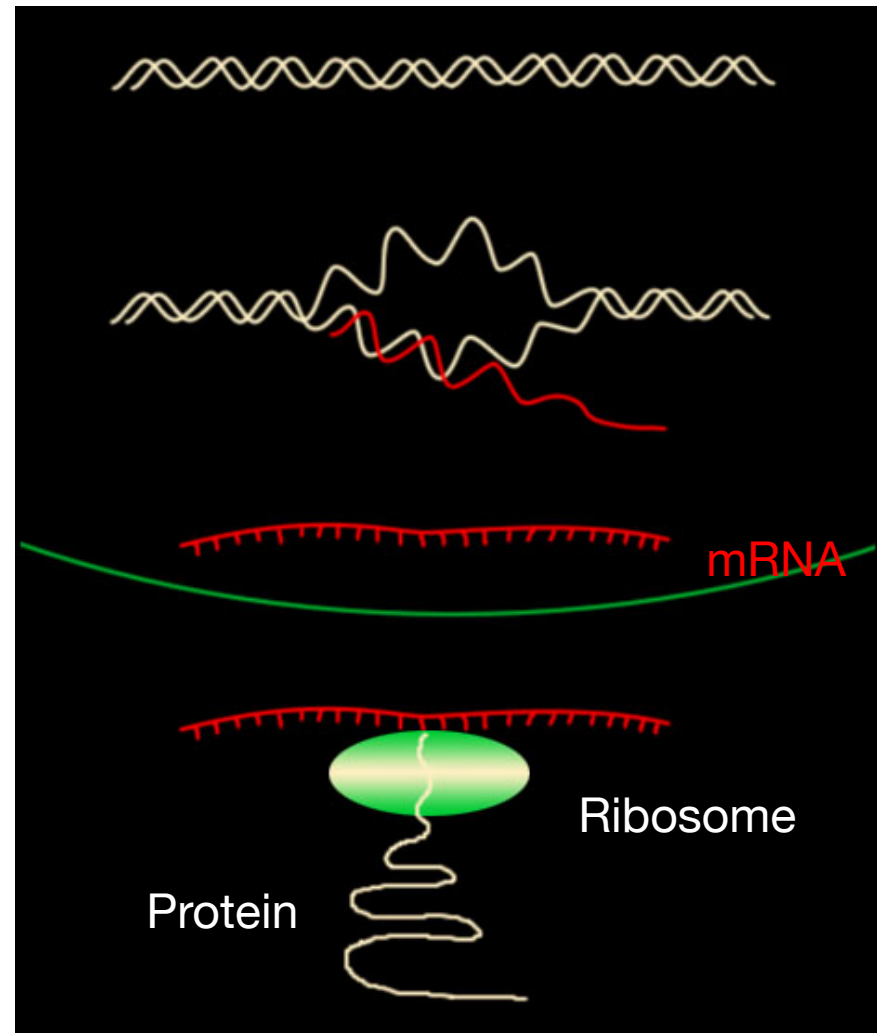
Phenotype



From Genes to Proteins

Transcription:
DNA to mRNA

Translation:
mRNA to Proteins



Proteins

Proteins are the “workhorses” of cells

- To understand how cells work is to understand proteins

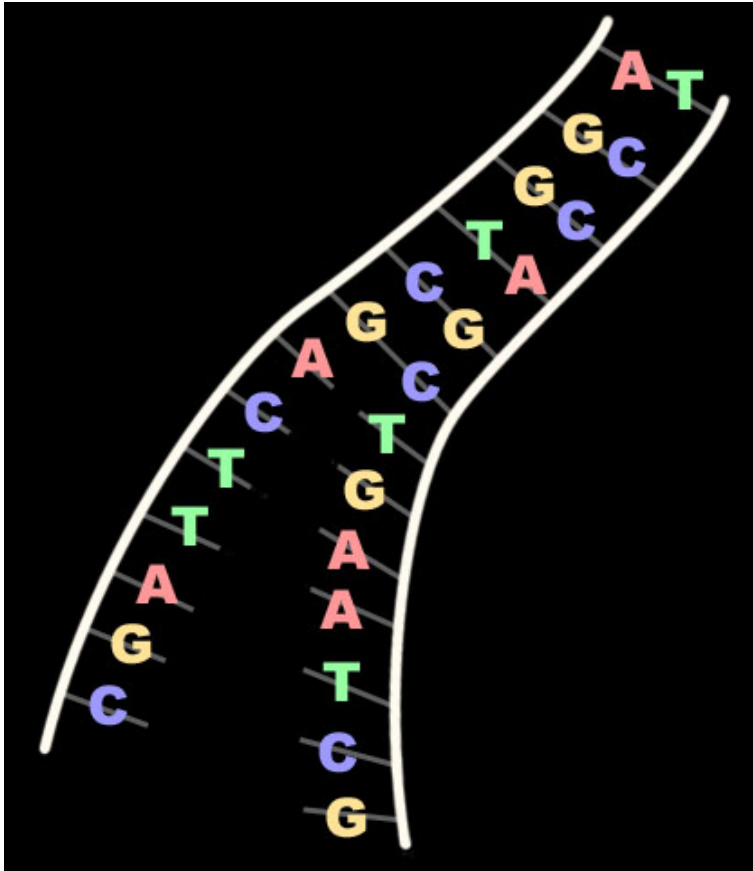
Understanding proteins and cells is key for finding disease treatments and cures

- Modern drug development is centered on affecting proteins (receptors, hormones, etc.)

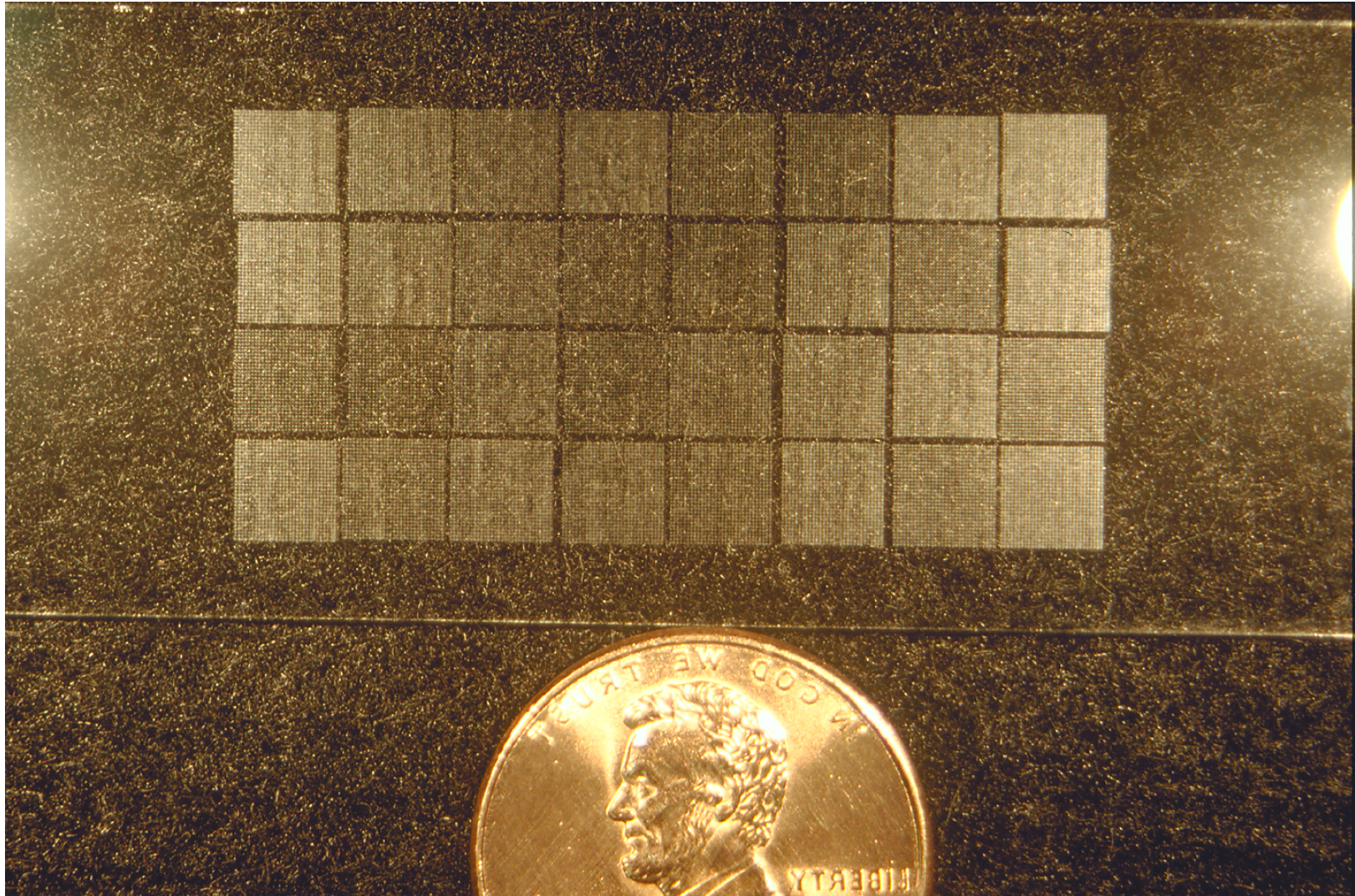
But... Proteins are hard to study directly, so microarrays look at the mRNA instead.

Hybridization

Expression microarrays use the fact that complementary strands will hybridize (attach) to each other



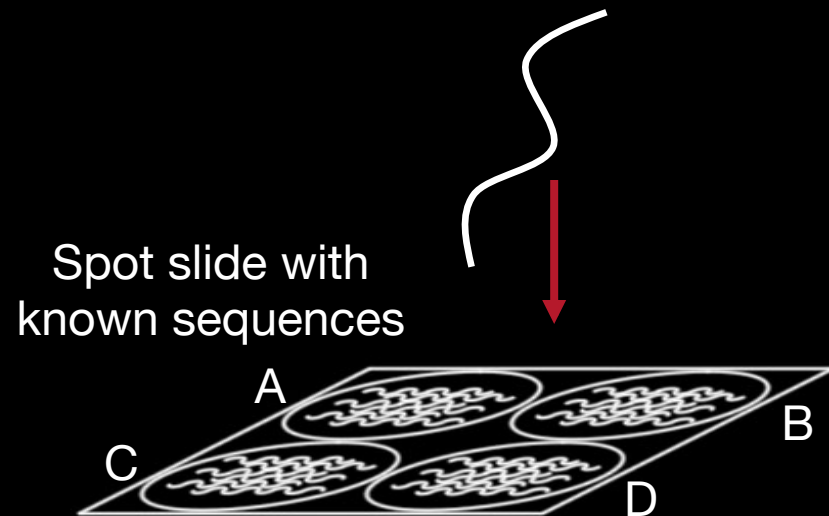
Early cDNA microarray (18,000 clones)



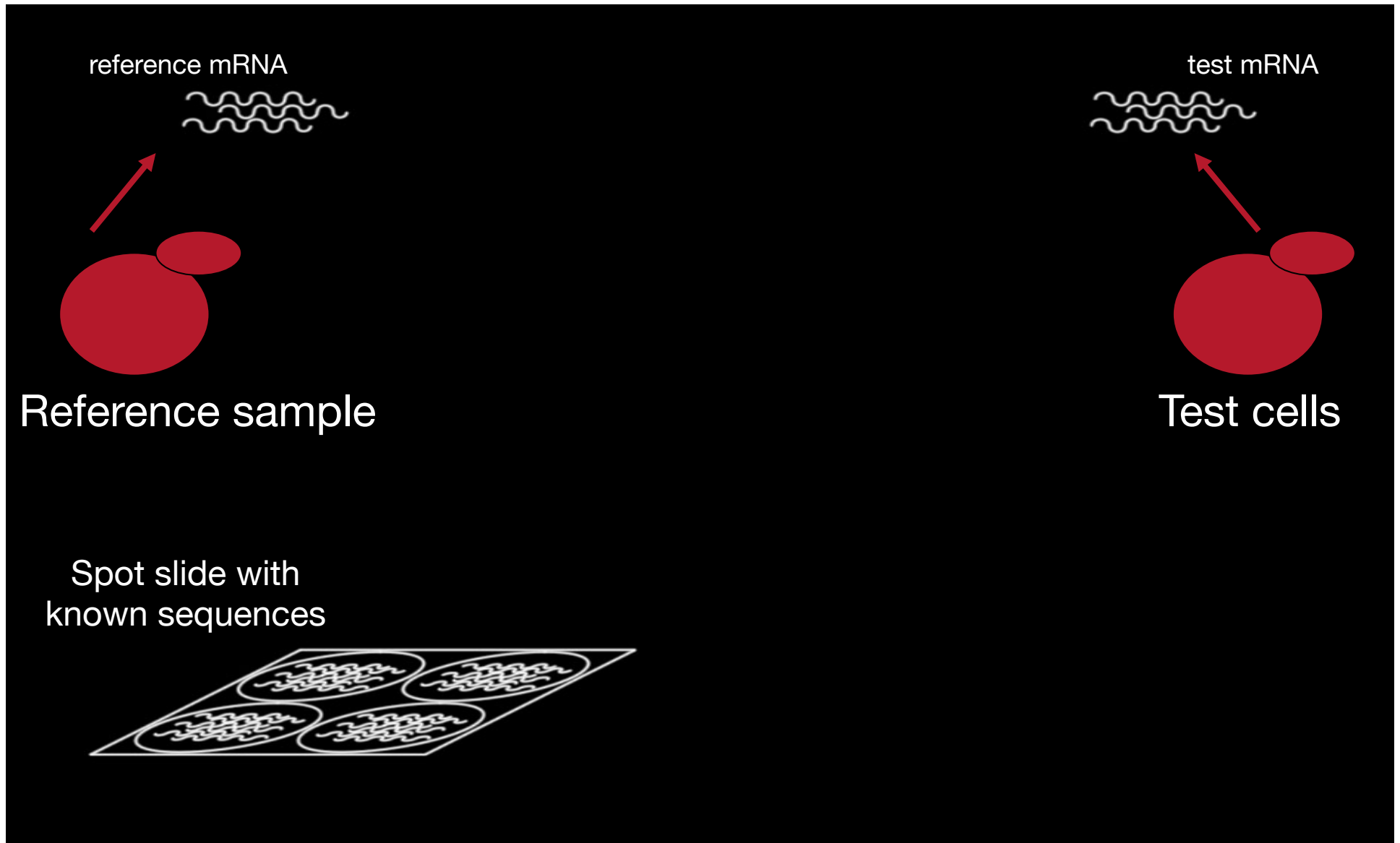
Microarray Methodology



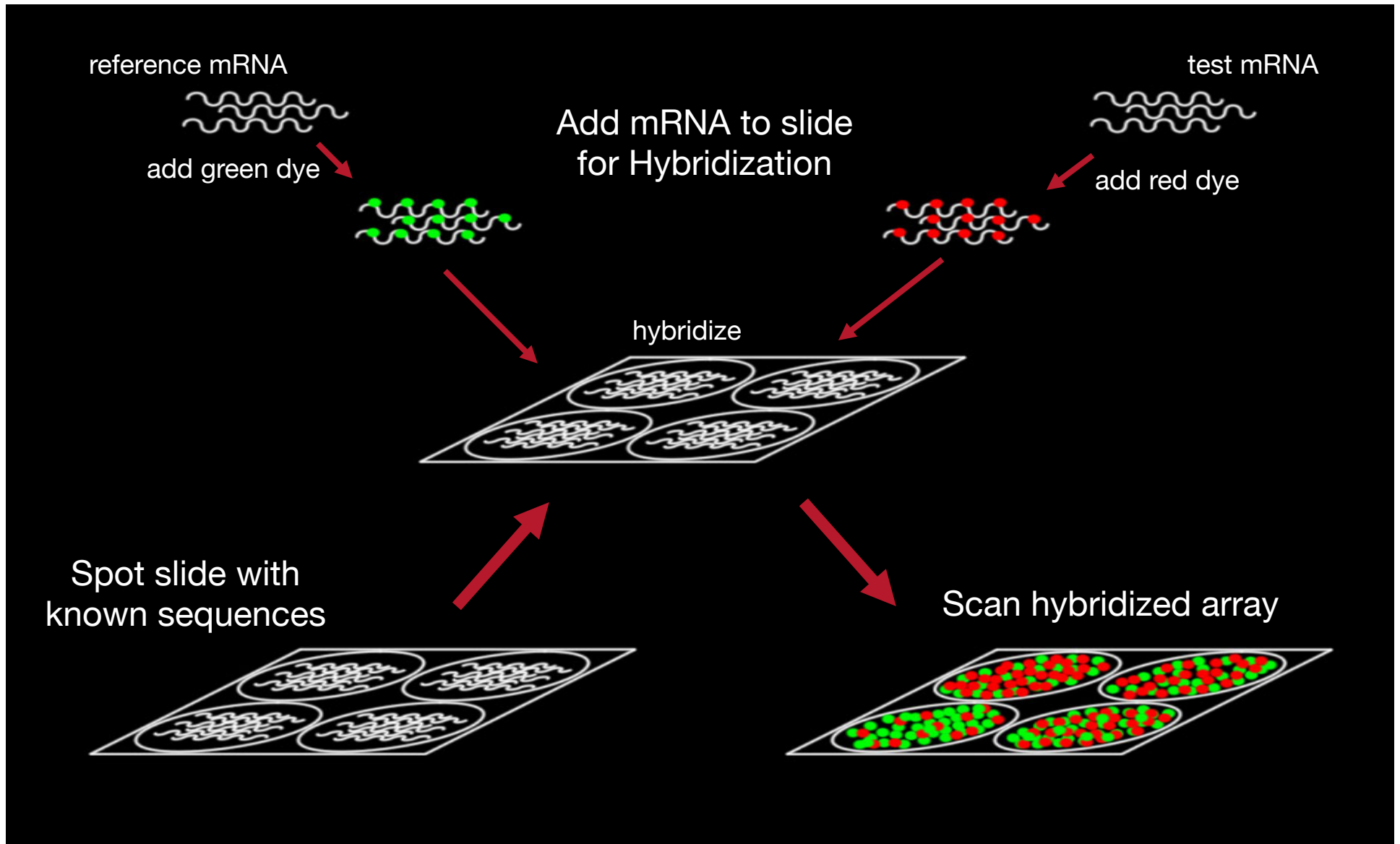
Microarray Methodology



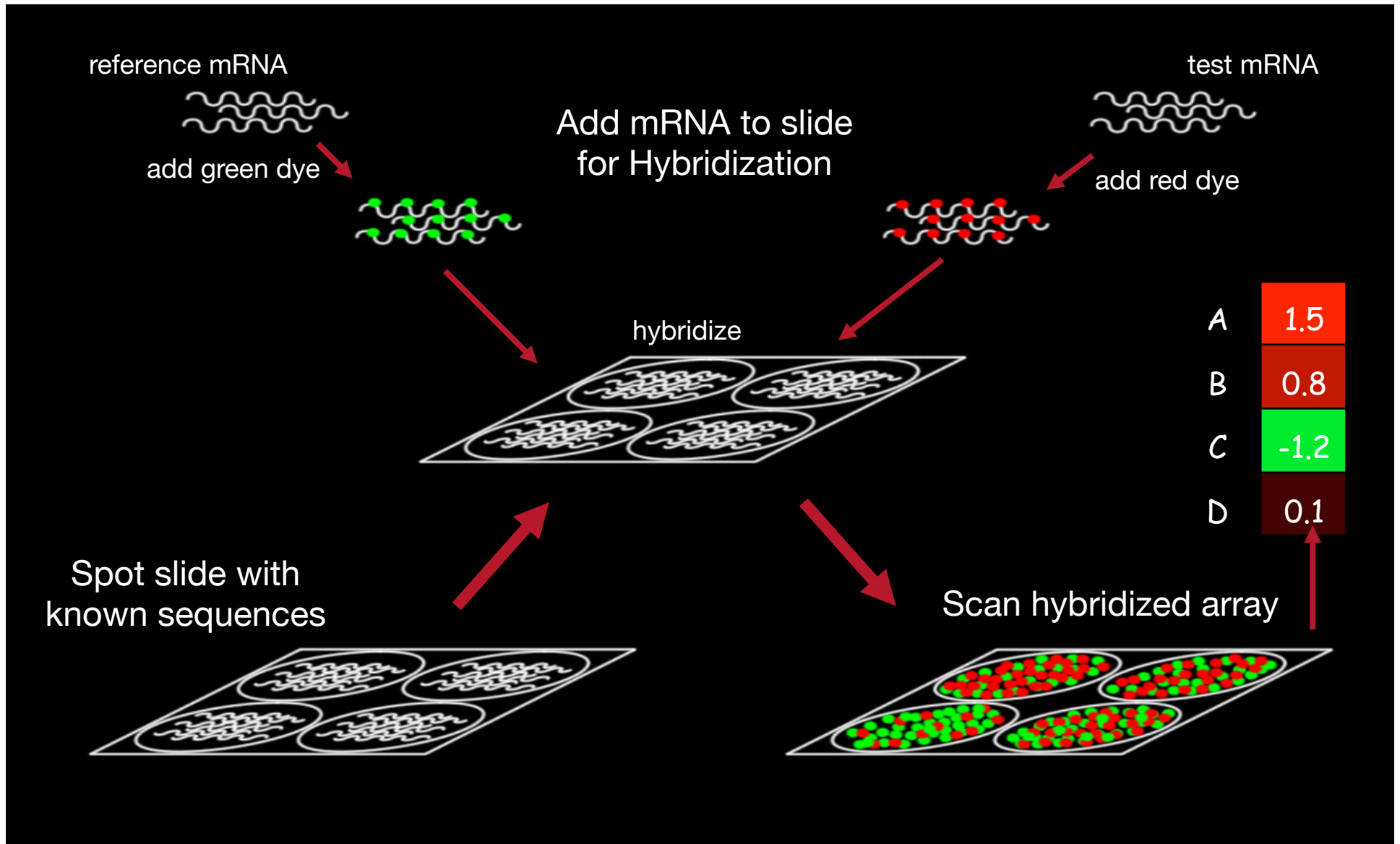
Microarray Methodology



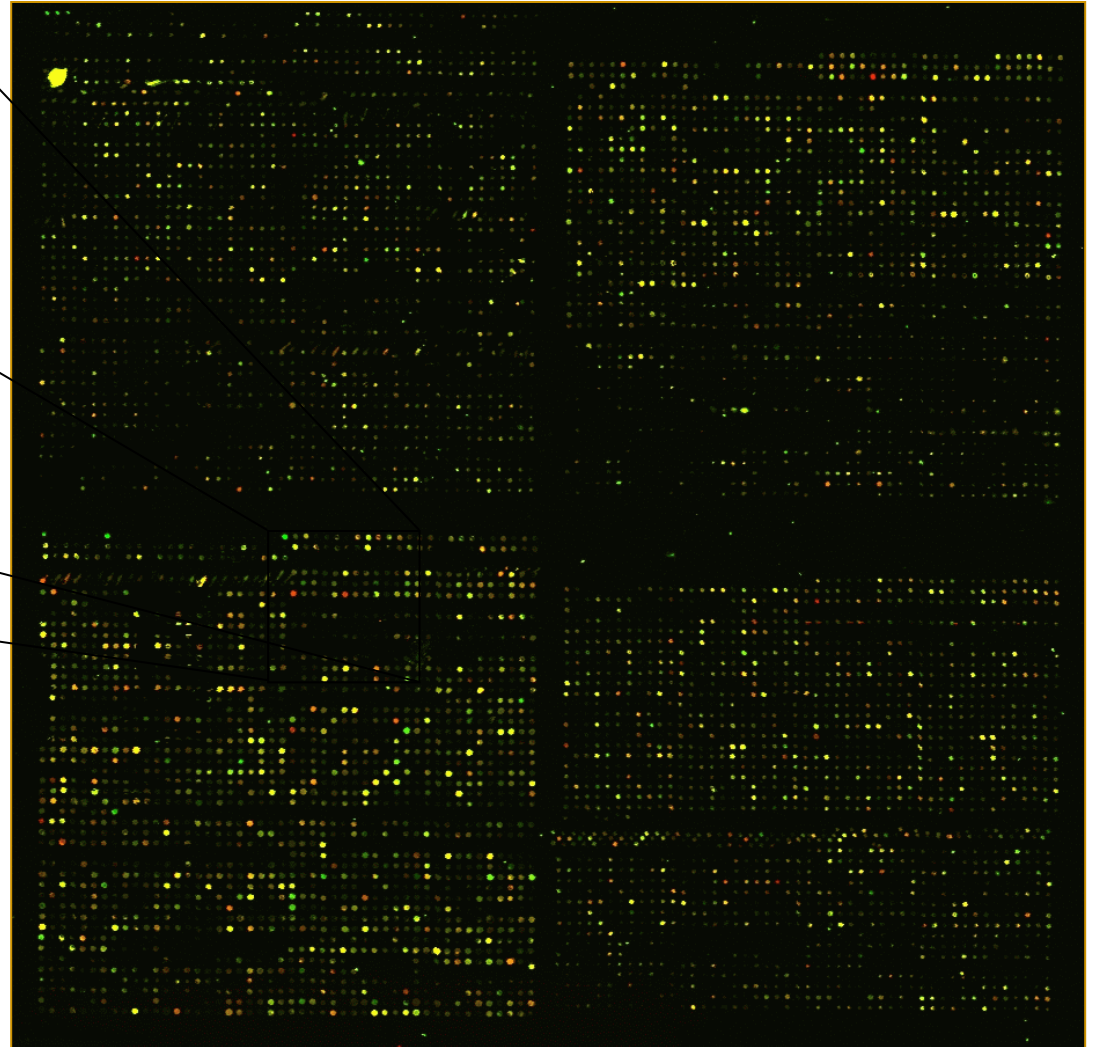
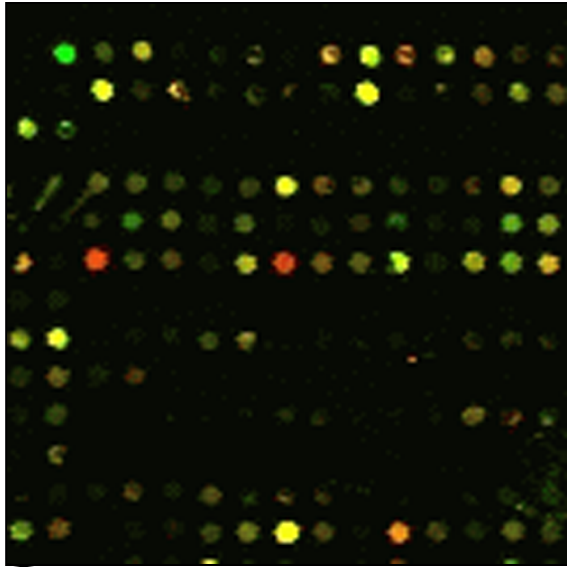
Microarray Methodology



Microarray Methodology



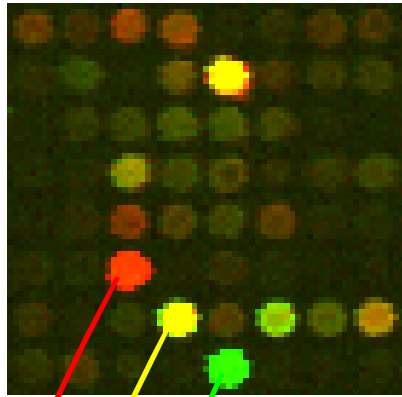
Microarray Outputs



Measure amounts of green and red dye on each spot

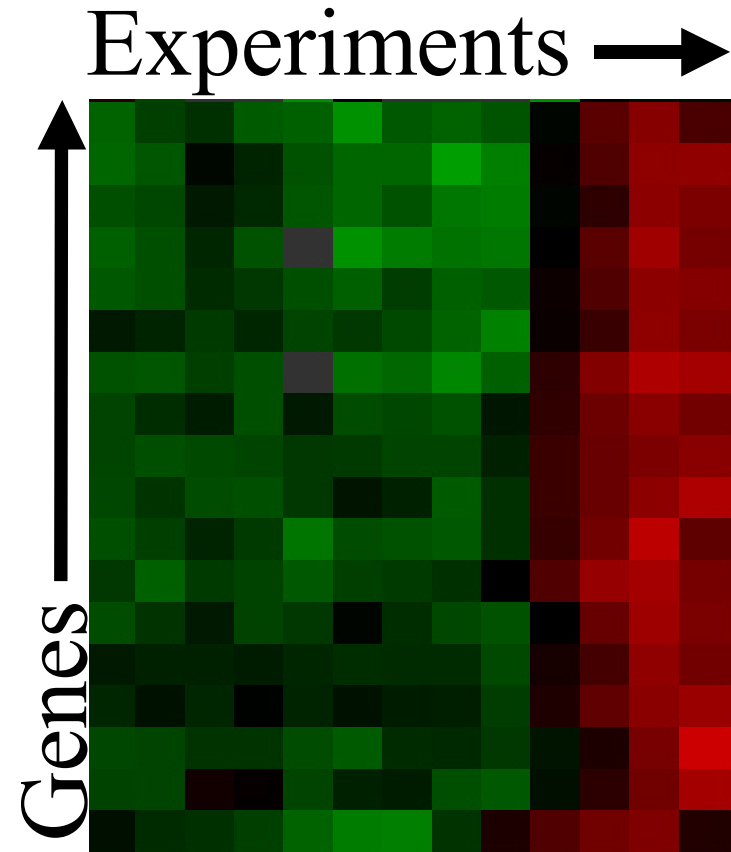
Represent level of expression as a log ratio between these amounts

Extracting Data



●	200	10000	50.00	5.64	■
●	4800	4800	1.00	0.00	■
●	9000	300	0.03	-4.91	■

Cy3 Cy5 $\frac{\text{Cy5}}{\text{Cy3}}$ $\log_2\left(\frac{\text{Cy5}}{\text{Cy3}}\right)$

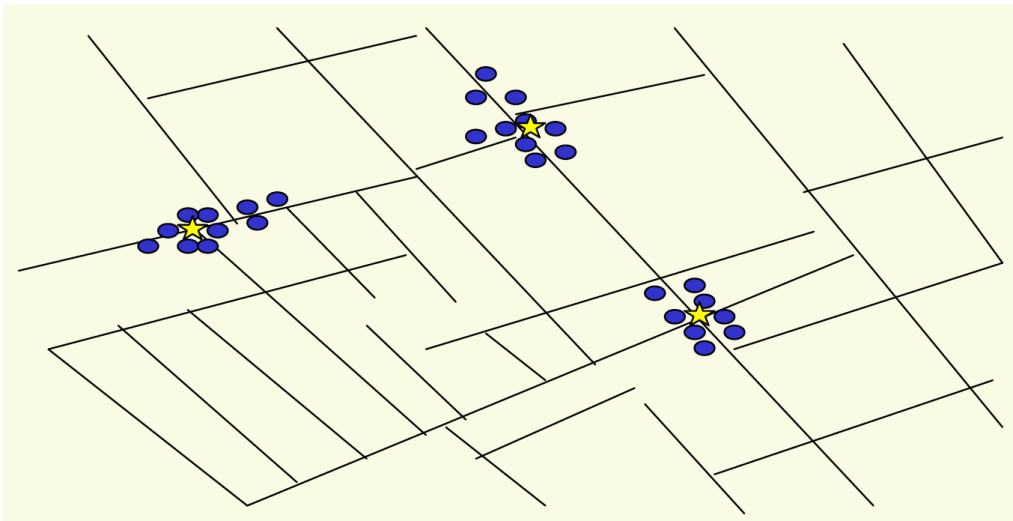


Some questions you can tackle with high-throughput gene-expression

Large-scale study of biological processes

- What is going on in the cell at a certain point in time?
 - what genes/pathways are active?
- On a genomic level, what accounts for differences between phenotypes?
 - which genes/pathways are activated in stress response?

Clustering

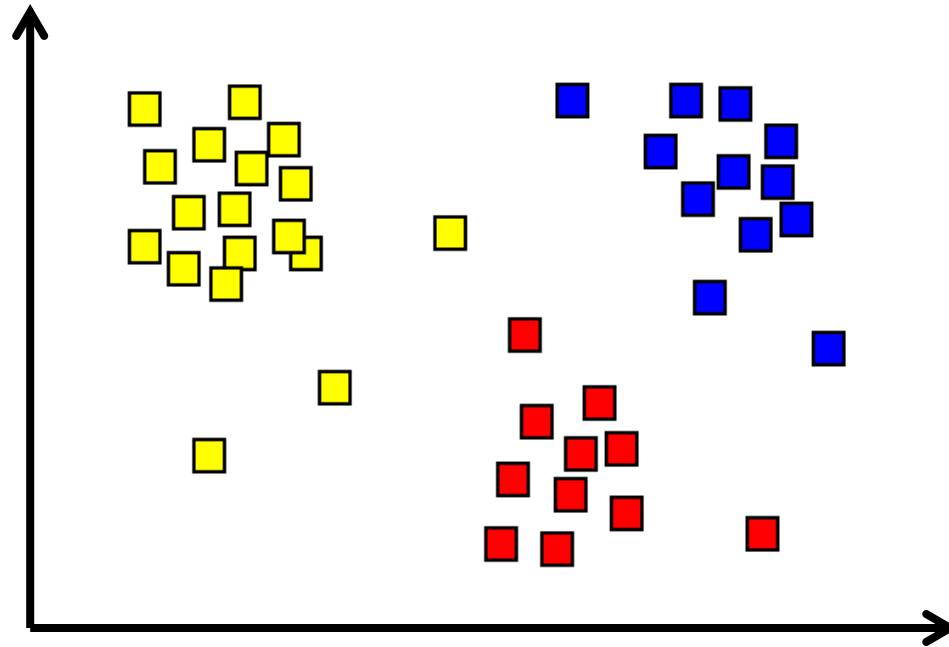


Outbreak of cholera deaths on map in 1850s.
Reference: Nina Mishra, HP Labs

History: London physicist John Snow plotted outbreak of cholera deaths on map in 1850s. Location indicated that clusters were around certain intersections with polluted wells; this exposed the problem and solution!

What is clustering?

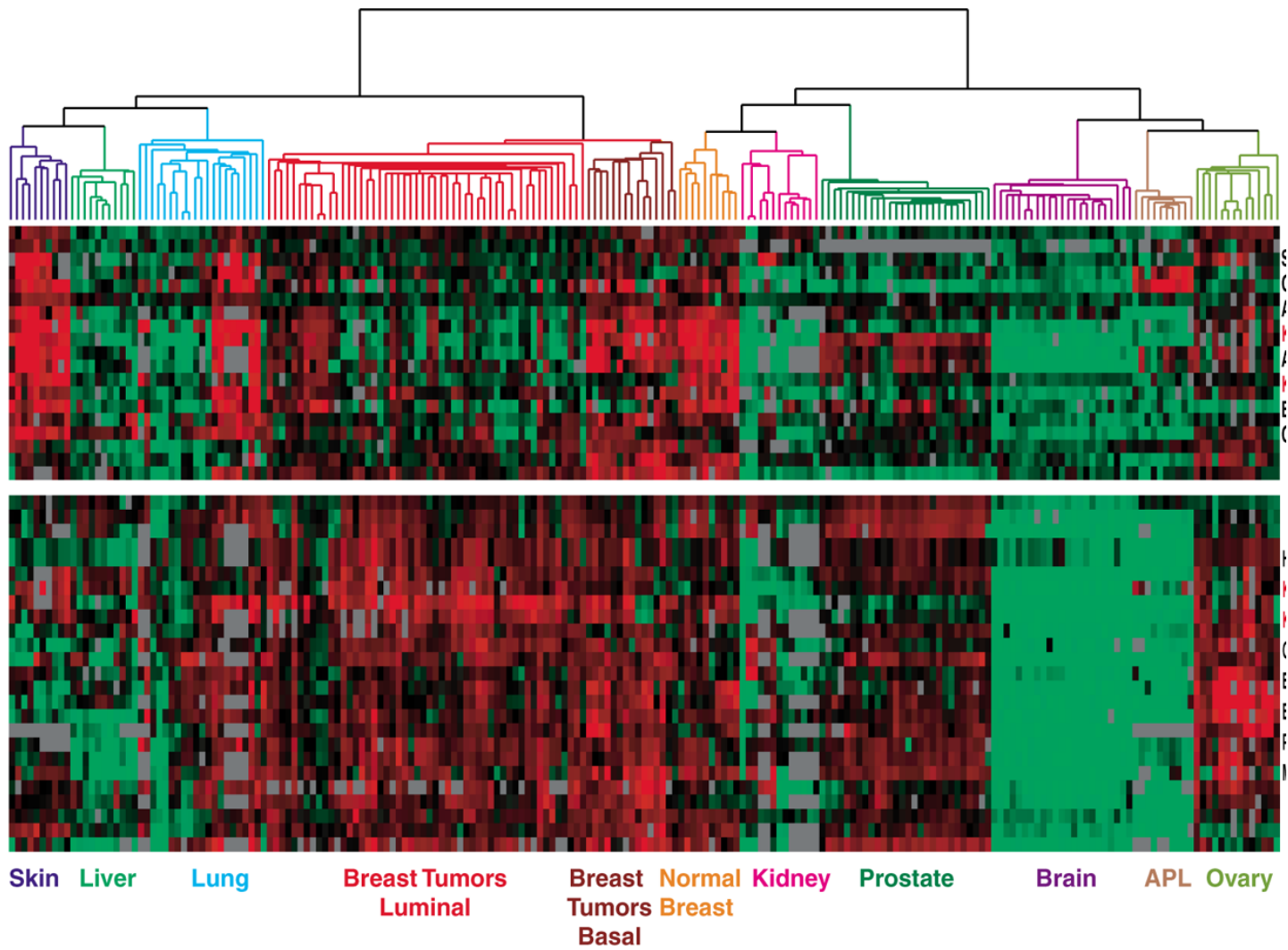
Reordering of vectors in a dataset so that similar patterns are next to each other



"Cluster-2" by Cluster-2.gif: hellispderivative work: Wgabrie (talk) - Cluster-2.gif. Licensed under Public Domain via Wikimedia Commons - <http://commons.wikimedia.org/wiki/File:Cluster-2.svg#mediaviewer/File:Cluster-2.svg>

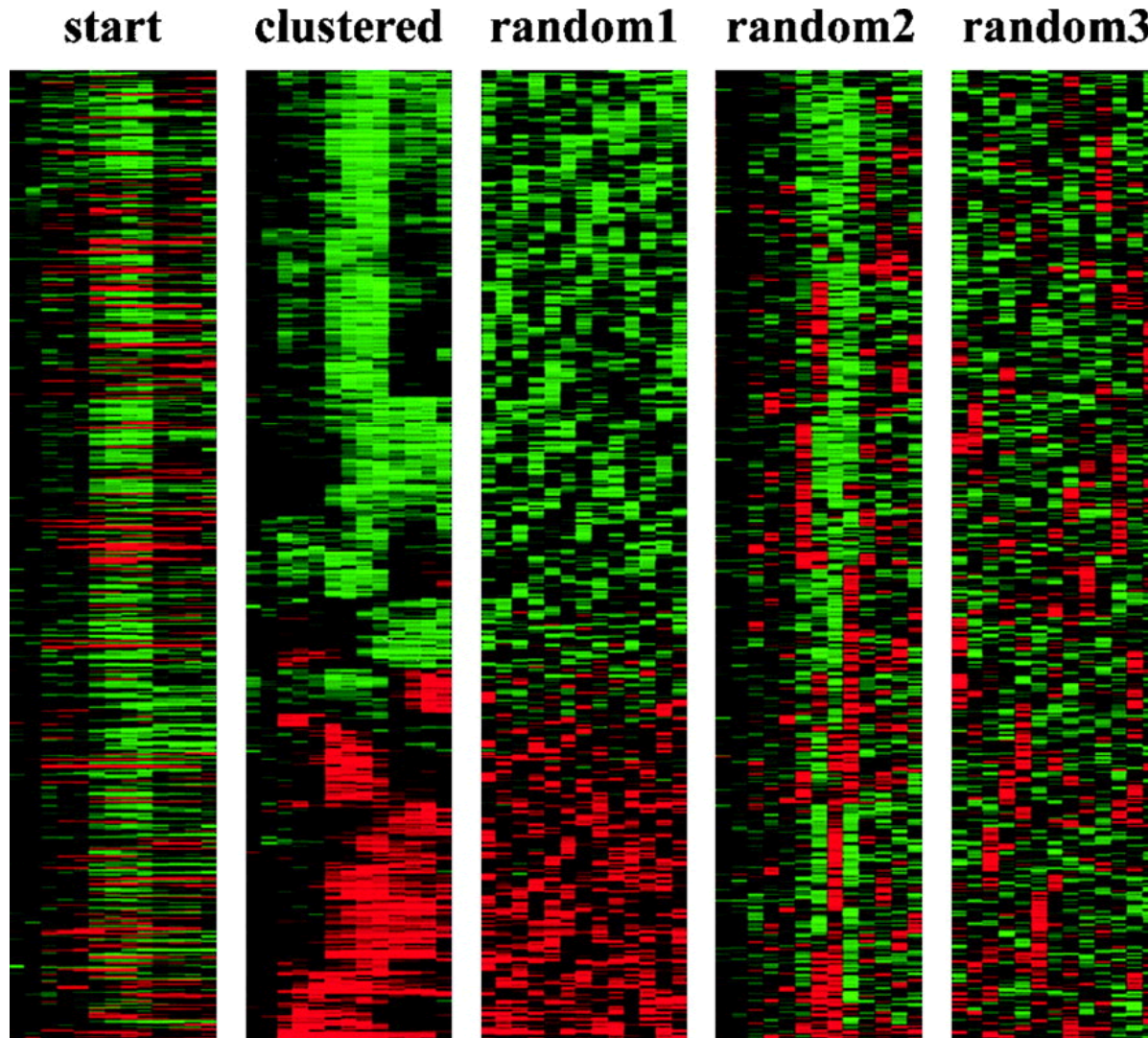
Why cluster microarray data?

- **Guilt-by-association:** if unknown gene i is similar in expression to known gene j , maybe they are involved in the same/related pathway
- **Dimensionality reduction:** datasets are too big to be able to get information out without reorganizing the data



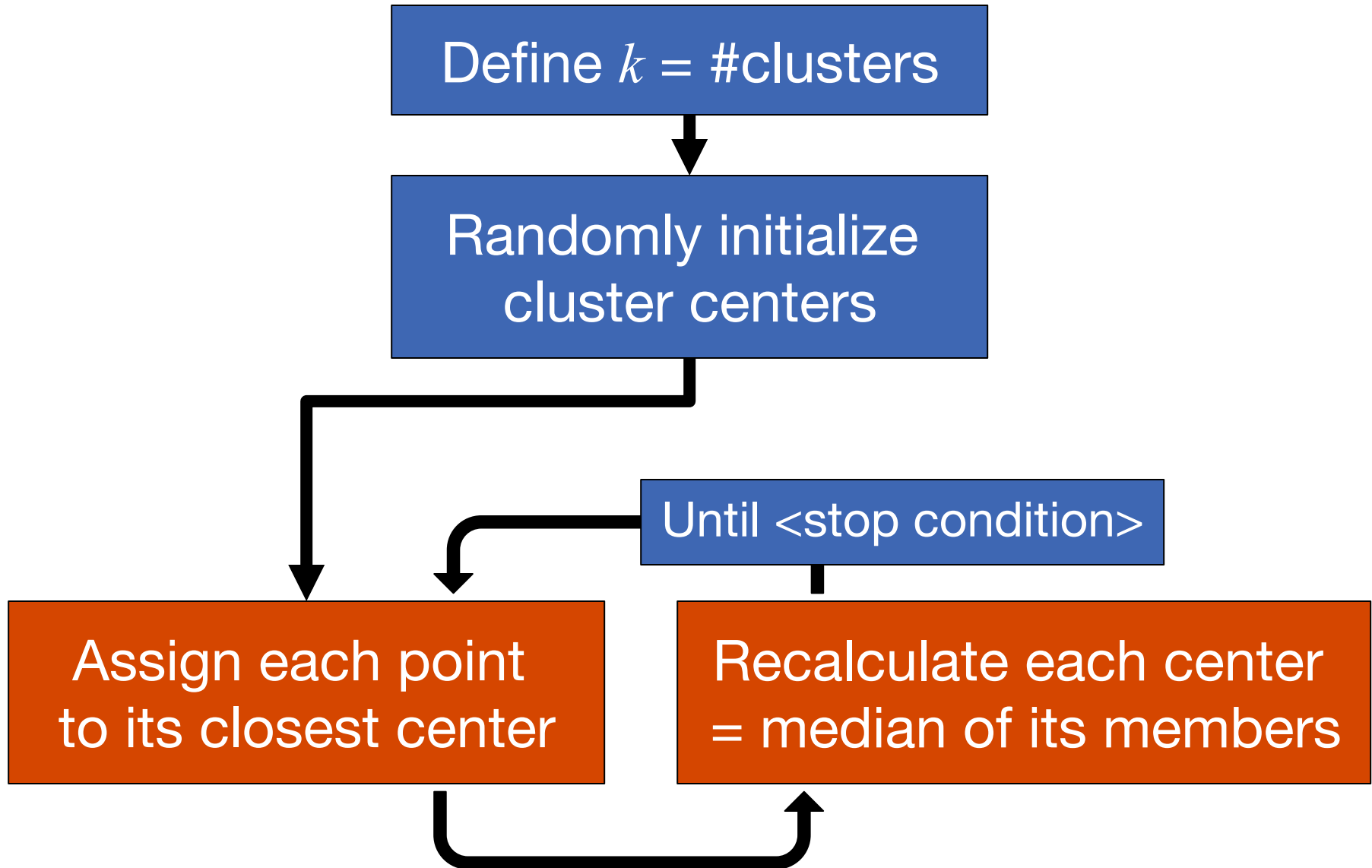
Botstein & Brown group

Clustering Random vs Biological Data



Challenge:
when is
clustering
“real”?

K-means clustering



K-means clustering

DEMO

<http://www.naftaliharris.com/blog/visualizing-k-means-clustering/>

K-means clustering

Conceptually similar to Expectation-Maximization

EM iteration alternates between 2 two steps:

1. E step: Creates a function for the expectation of the log-likelihood evaluated using the current estimate for the parameters, and
2. M step: Computes parameters maximizing the expected log-likelihood found on the E step.

These parameter-estimates are then used to determine the distribution of the latent variables in the next E step.

K-means clustering

Stopping condition

- Until the change in centers is less than $\langle \text{constant} \rangle$
- Until all genes get assigned to the same partition twice in a row
- Until some minimal number of genes (e.g. 90%) get assigned to the same partition twice in a row

K-means clustering

Some issues

- Have to set k ahead of time
- Prefers clusters of approx. similar sizes
- Each gene only belongs to 1 cluster
- Genes assigned to clusters on the basis of all experiments

Hierarchical clustering

- Imposes hierarchical structure on all of the data
- Easy visualization of similarities and differences between genes (experiments) and clusters of genes (experiments)

Hierarchical clustering

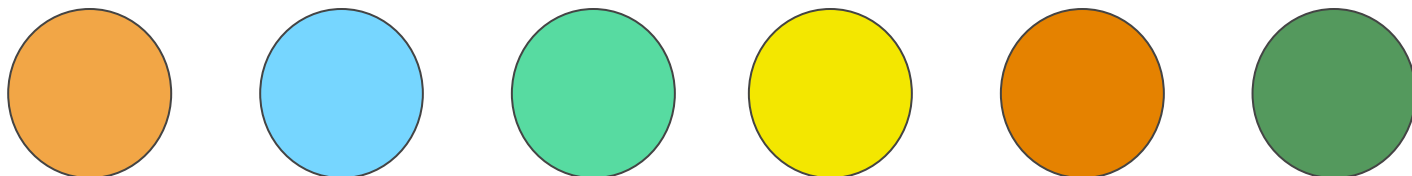
Start with each pattern
in its own cluster

Until all patterns
are merged into a
single cluster

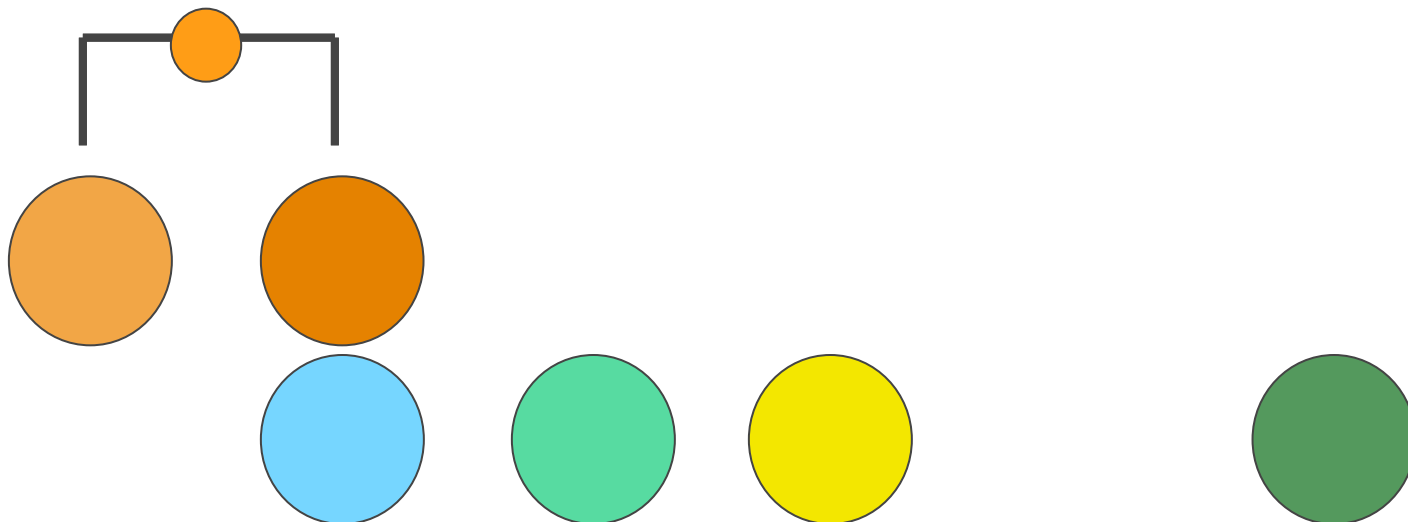
Join patterns that are
most similar

Compare joined patterns
to all un-joined patterns

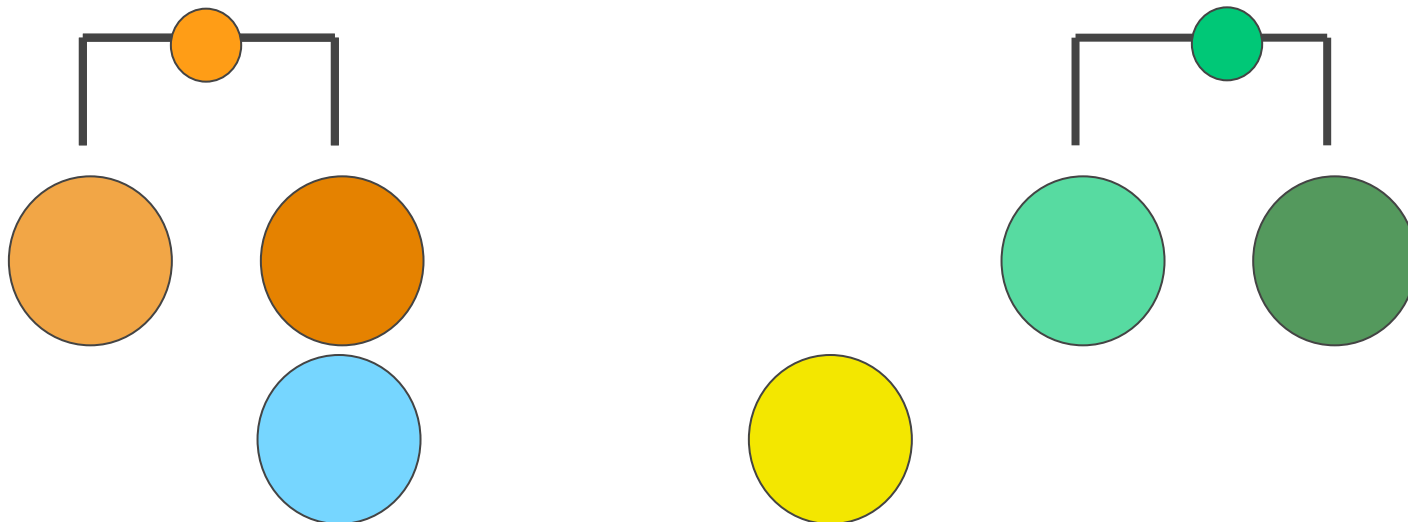
Hierarchical clustering



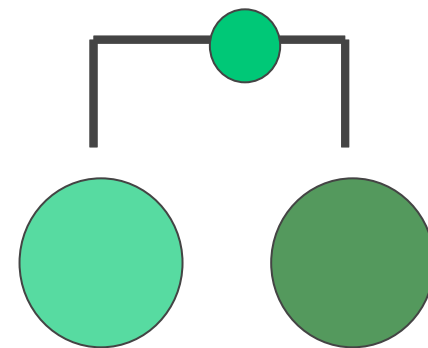
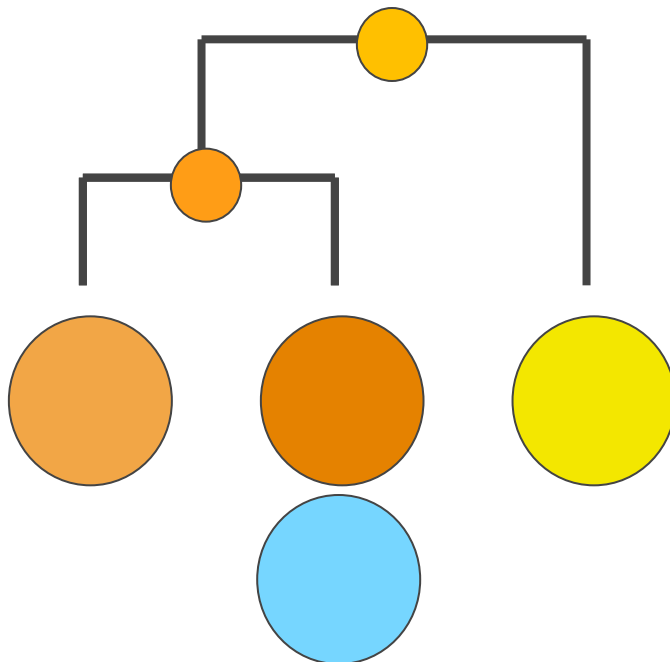
Hierarchical clustering



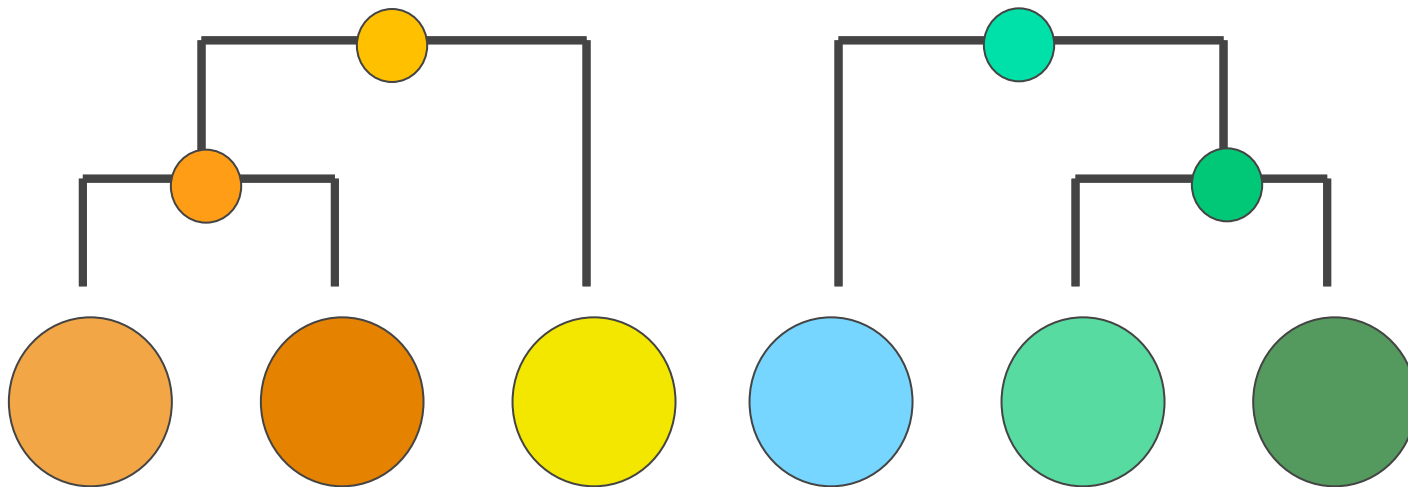
Hierarchical clustering



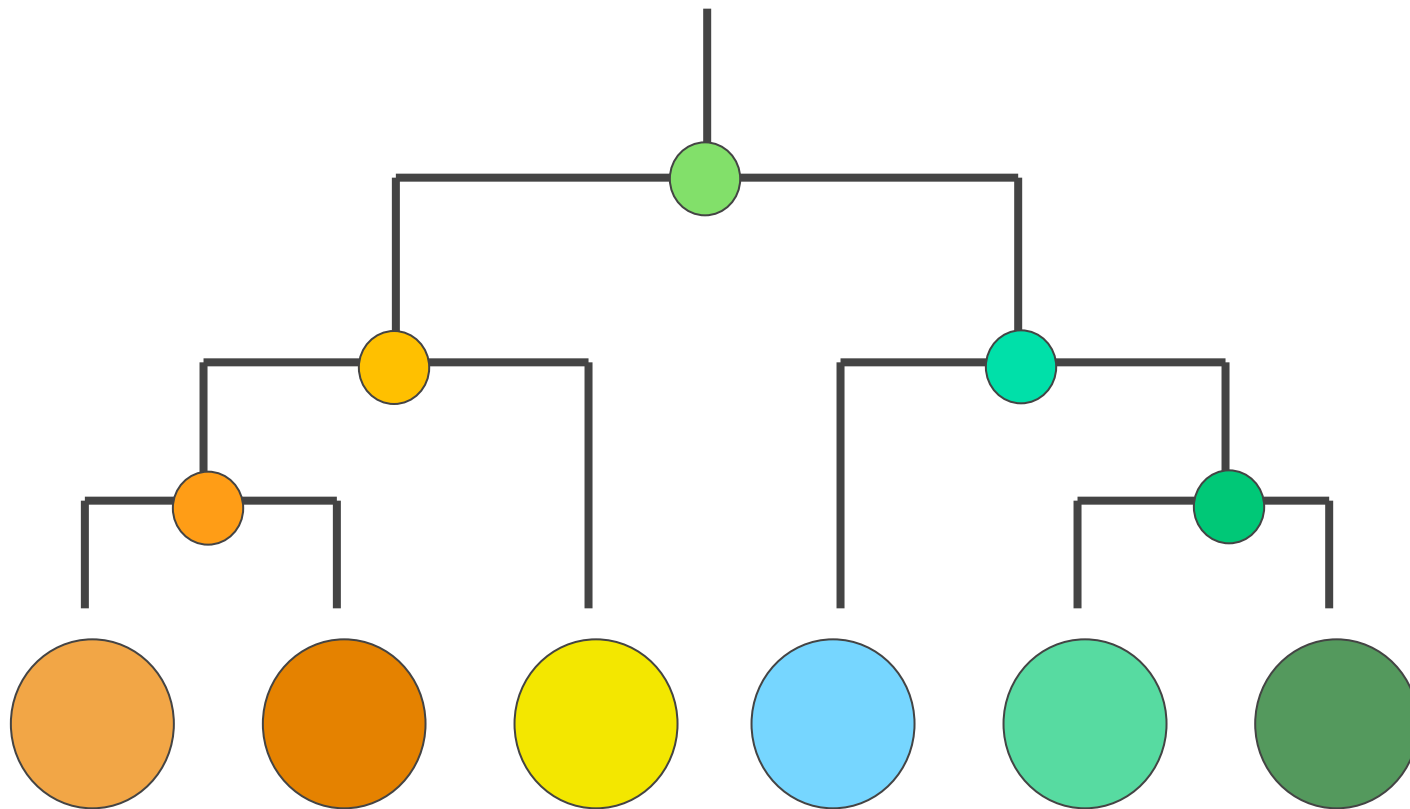
Hierarchical clustering



Hierarchical clustering

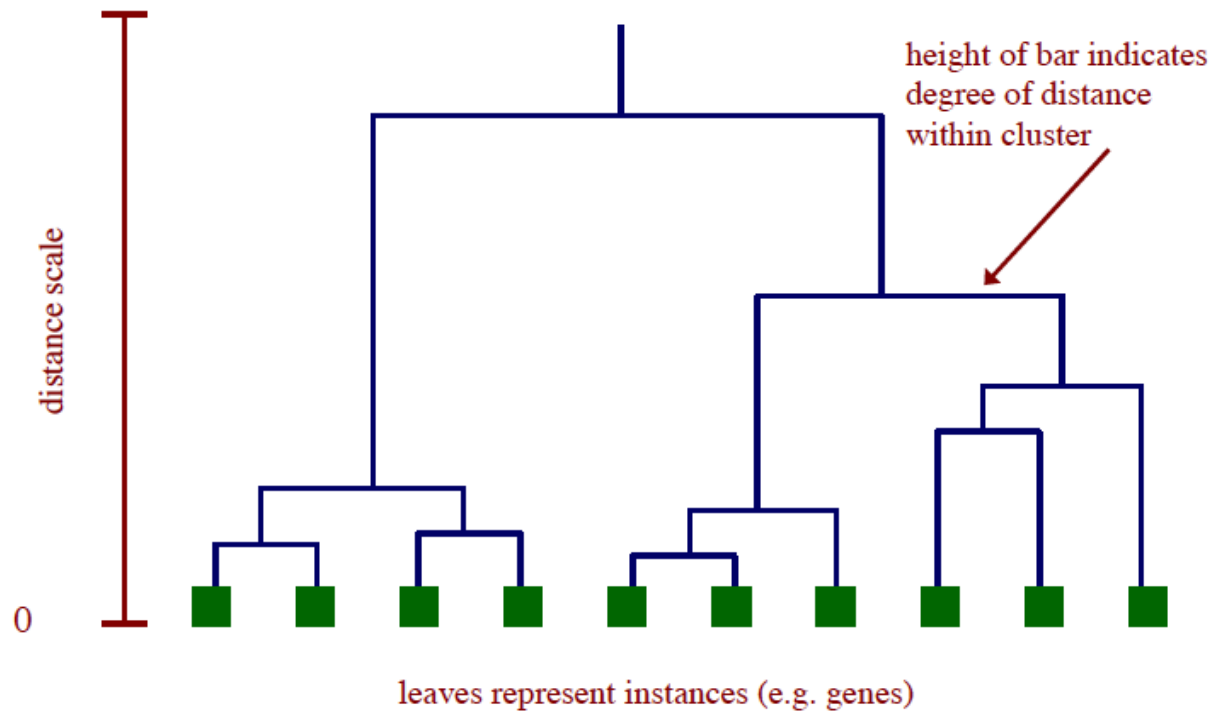


Hierarchical clustering



Dendrogram

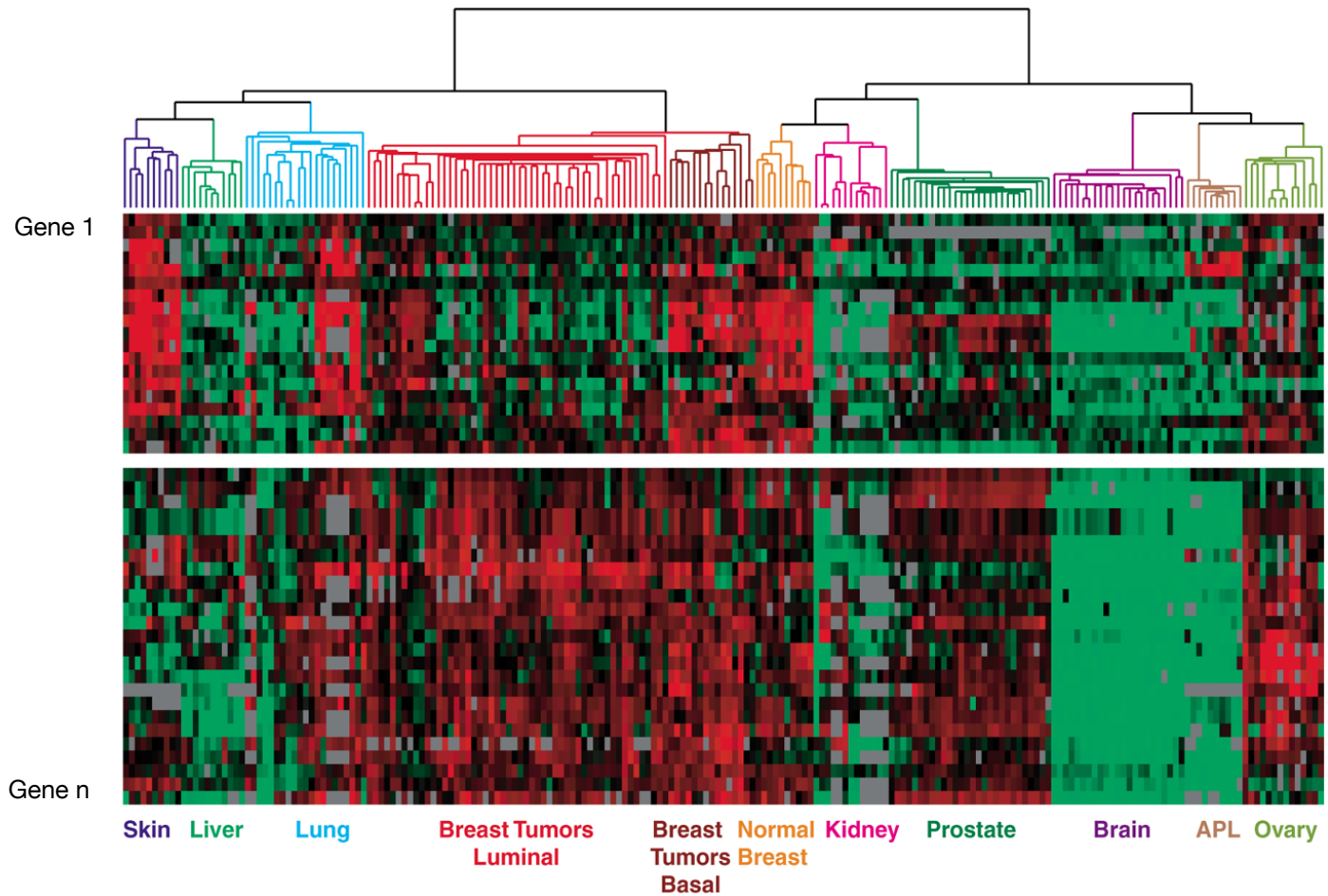
- Leaves = genes.
- Internal nodes = hypothetical ancestors.



Reference: <http://www.biostat.wisc.edu/bmi576/fall-2003/lecture13.pdf>

Dendrogram of Human tumors

Tumors in similar tissues cluster together.



Reference: Botstein & Brown group

■ gene over expressed
■ gene under expressed