

CAP 5510: Introduction to Bioinformatics  
CGS 5166: Bioinformatics Tools

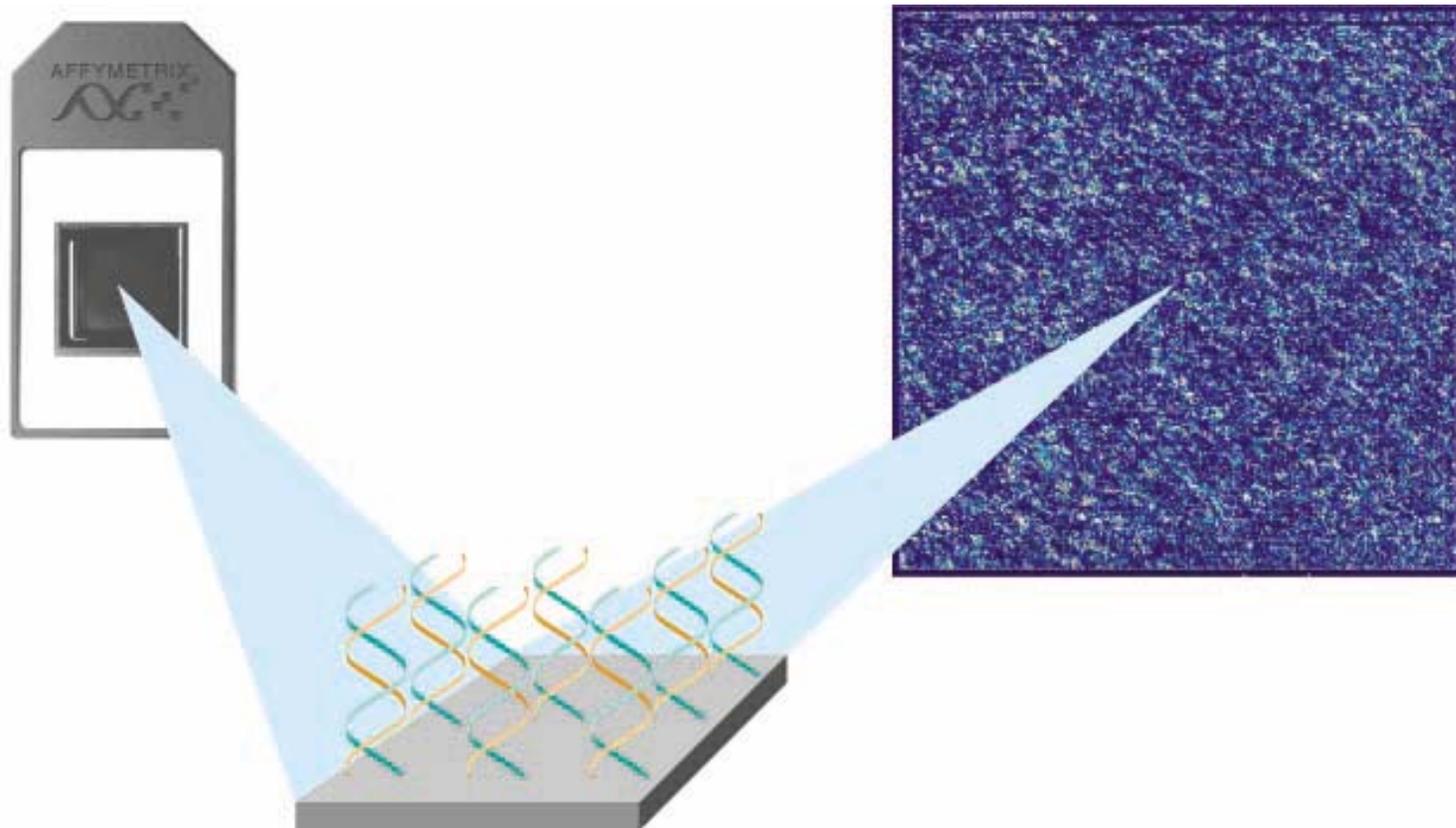
**Giri Narasimhan**

ECS 254; Phone: x3748

[giri@cis.fiu.edu](mailto:giri@cis.fiu.edu)

[www.cis.fiu.edu/~giri/teach/BioinfS11.html](http://www.cis.fiu.edu/~giri/teach/BioinfS11.html)

# DNA Chips & Images



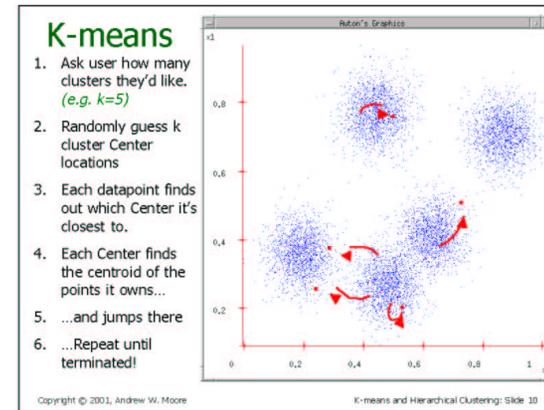
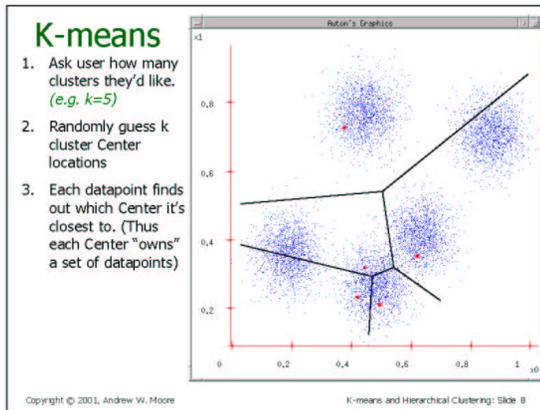
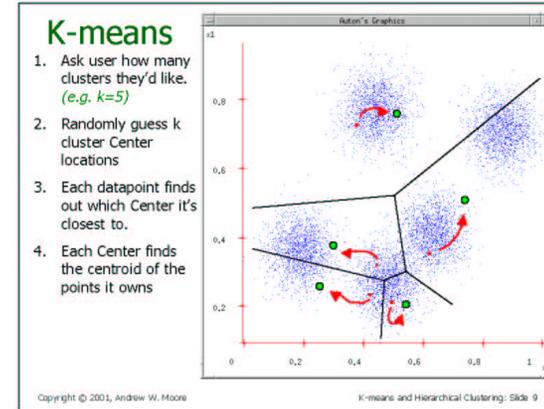
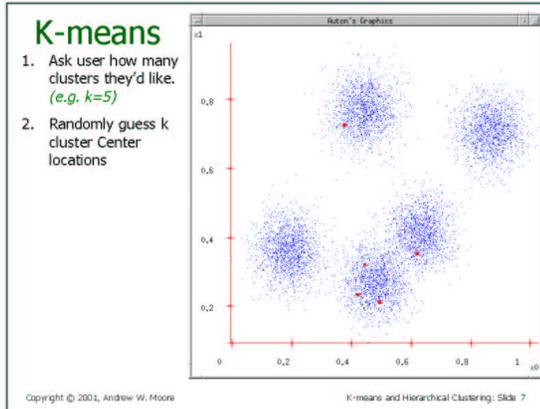
# Microarray Data

<i>Gene</i>	<i>Expression Level</i>
<i>Gene1</i>	
<i>Gene2</i>	
<i>Gene3</i>	
...	

# K-Means Clustering: Example

Example from Andrew Moore's tutorial on Clustering.

Start



4

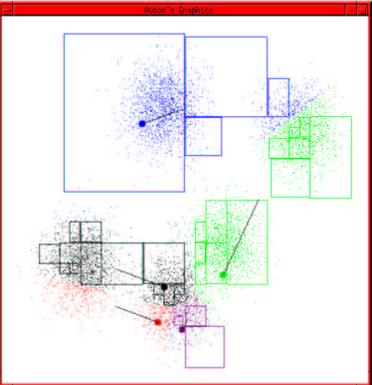
5

## K-means Start

Advance apologies: in Black and White this example will deteriorate

Example generated by Dan Pelleg's super-duper fast K-means system:

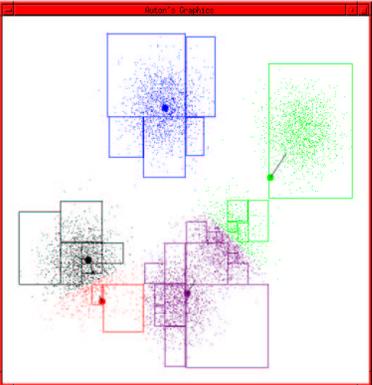
*Dan Pelleg and Andrew Moore. Accelerating Exact k-means Algorithms with Geometric Reasoning. Proc. Conference on Knowledge Discovery in Databases 1999, (KDD99) (available on [www.autonlab.org/pap.html](http://www.autonlab.org/pap.html))*



Copyright © 2001, Andrew W. Moore K-means and Hierarchical Clustering: Slide 11

## K-means continues

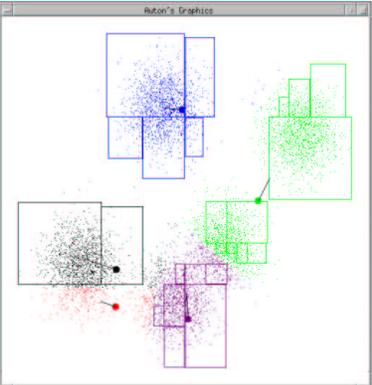
...



Copyright © 2001, Andrew W. Moore K-means and Hierarchical Clustering: Slide 13

## K-means continues

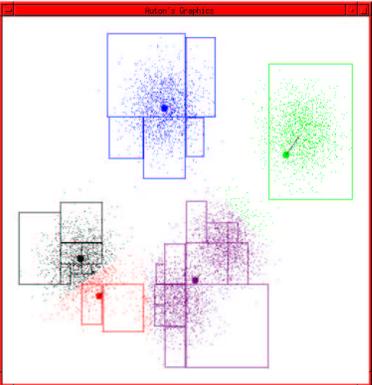
...



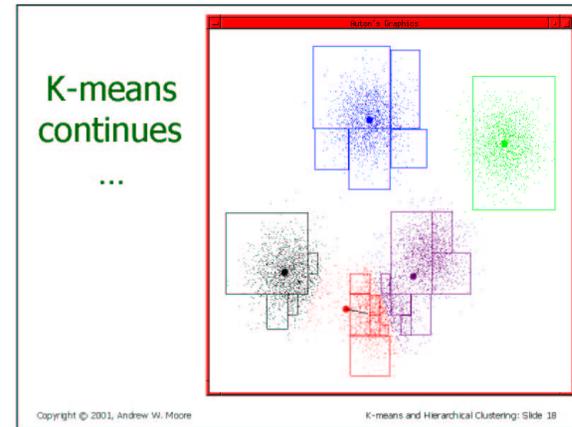
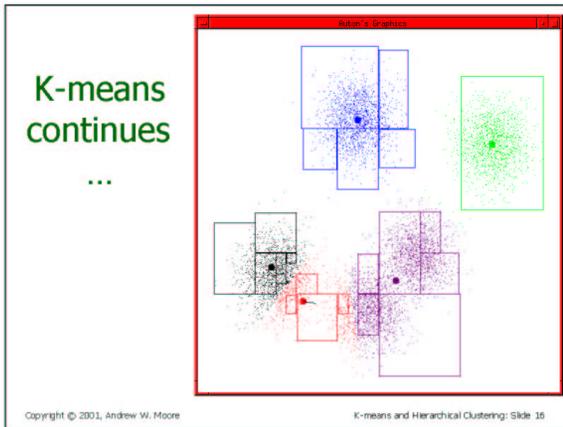
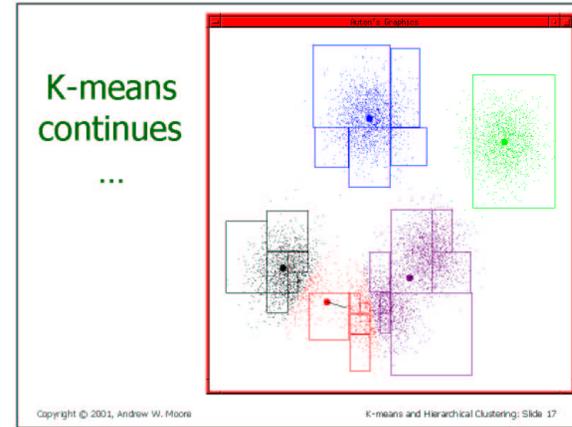
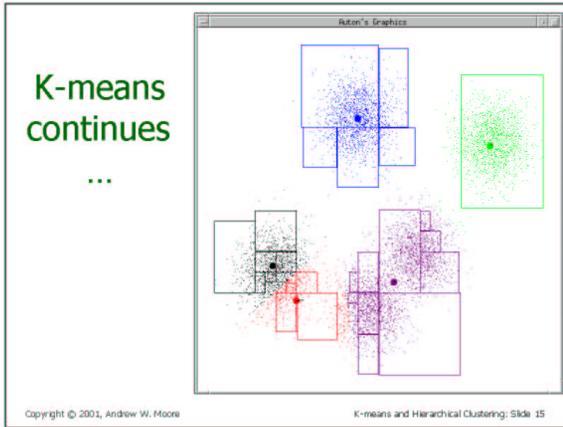
Copyright © 2001, Andrew W. Moore K-means and Hierarchical Clustering: Slide 12

## K-means continues

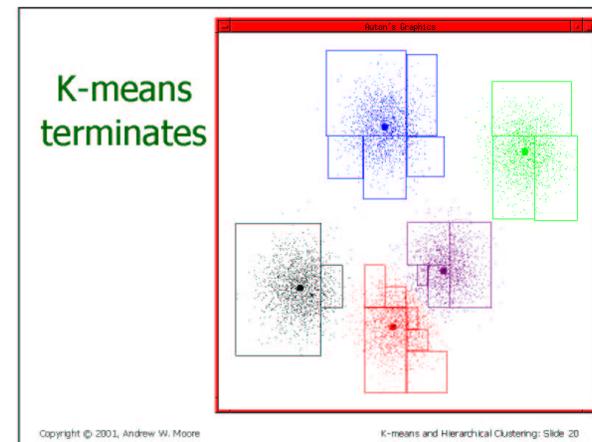
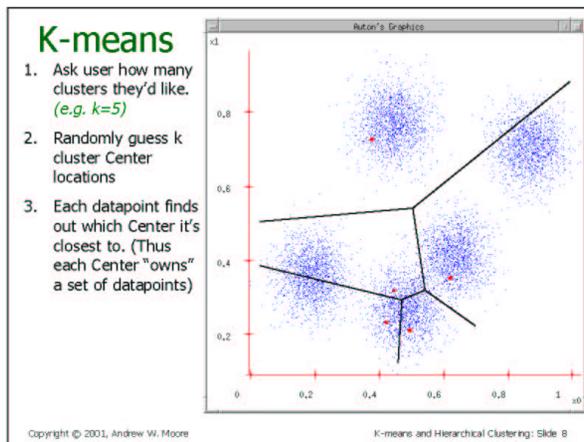
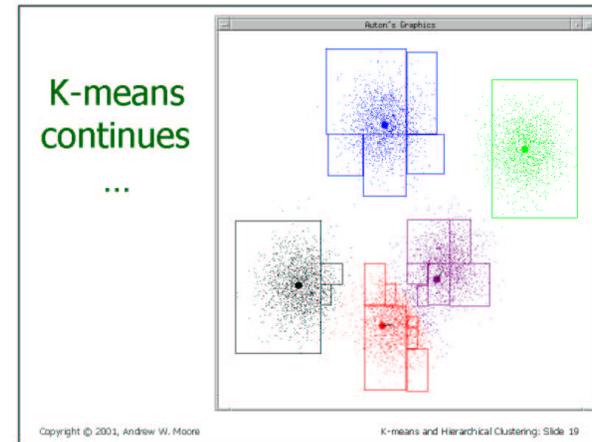
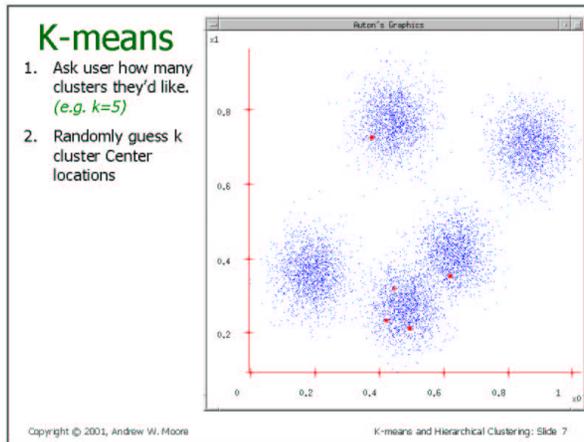
...



Copyright © 2001, Andrew W. Moore K-means and Hierarchical Clustering: Slide 14



Start



End

# K-Means Clustering [McQueen '67]

Repeat

- Start with randomly chosen cluster centers
- Assign points to give greatest increase in score
- Recompute cluster centers
- Reassign points

until (no changes)

Try the applet at: [http://home.dei.polimi.it/matteucc/Clustering/tutorial\\_html/AppletH.html](http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/AppletH.html)

# Comparisons

## □ Hierarchical clustering

- Number of clusters not preset.
- Complete hierarchy of clusters
- Not very robust, not very efficient.

## □ K-Means

- Need definition of a **mean**. Categorical data?
- More efficient and often finds optimum clustering.

## Functionally related genes behave similarly across experiments

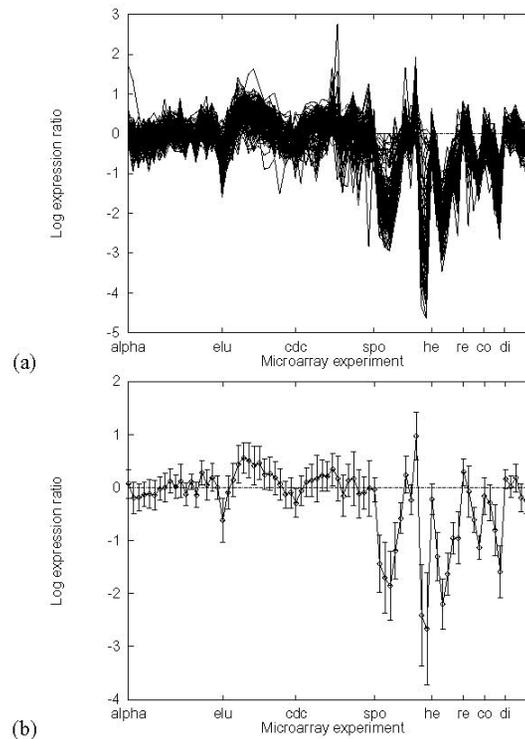


Figure 1: **Expression profiles of the cytoplasmic ribosomal proteins.** Figure (a) shows the expression profiles from the data in [Eisen et al., 1998] of 121 cytoplasmic ribosomal proteins, as classified by MYGD [MYGD, 1999]. The logarithm of the expression ratio is plotted as a function of DNA microarray experiment. Ticks along the X-axis represent the beginnings of experimental series. They are, from left to right, cell division cycle after synchronization with  $\alpha$  factor arrest (alpha), cell division cycle after synchronization by centrifugal elutriation (elu), cell division cycle measured using a temperature sensitive *cdc15* mutant (cdc), sporulation (spo), heat shock (he), reducing shock (re), cold shock (co), and diauxic shift (di). Sporulation is the generation of a yeast spore by meiosis. Diauxic shift is the shift from anaerobic (fermentation) to aerobic (respiration) metabolism. The medium starts rich in glucose, and yeast cells ferment, producing ethanol. When the glucose is used up, they switch to ethanol as a source for carbon. Heat, cold, and reducing shock are various ways to stress the yeast cell. Figure (b) shows the average, plus or minus one standard deviation, of the data in Figure (a).

# Self-Organizing Maps [Kohonen]

- ❑ Kind of neural network.
- ❑ Clusters data and find complex relationships between clusters.
- ❑ Helps reduce the dimensionality of the data.
- ❑ Map of 1 or 2 dimensions produced.
- ❑ Unsupervised Clustering
- ❑ Like K-Means, except for visualization

# SOM Architectures

- 2-D Grid
- 3-D Grid
- Hexagonal Grid

# SOM Algorithm

- Select SOM architecture, and initialize weight vectors and other parameters.
- **While** (stopping condition not satisfied) **do** for each input point  $x$ 
  - winning node  $q$  has weight vector **closest** to  $x$ .
  - **Update** weight vector of  $q$  and its **neighbors**.
  - **Reduce neighborhood size** and **learning rate**.

# SOM Algorithm Details

□ Distance between  $x$  and weight vector:  $\|x - w_i\|$

□ Winning node:  $q(x) = \min_i \|x - w_i\|$

□ Weight update function (for neighbors):

$$w_i(k+1) = w_i(k) + \mu(k, x, i)[x(k) - w_i(k)]$$

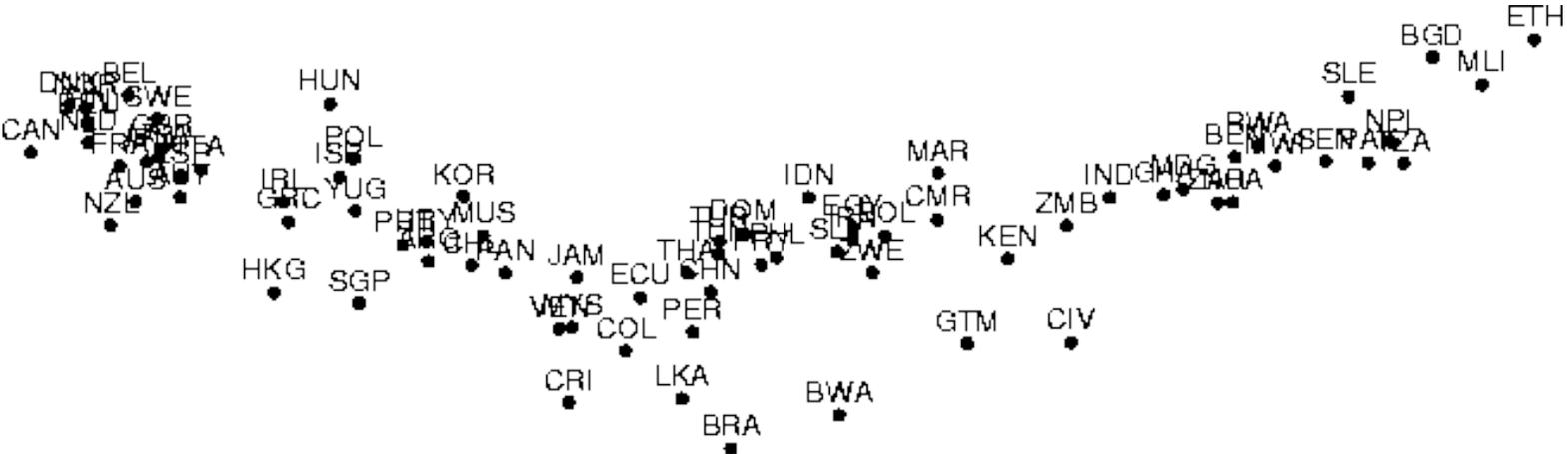
□ Learning rate:

$$\mu(k, x, i) = \eta_0(k) \exp\left(\frac{-\|r_i - r_{q(x)}\|^2}{\sigma^2}\right)$$

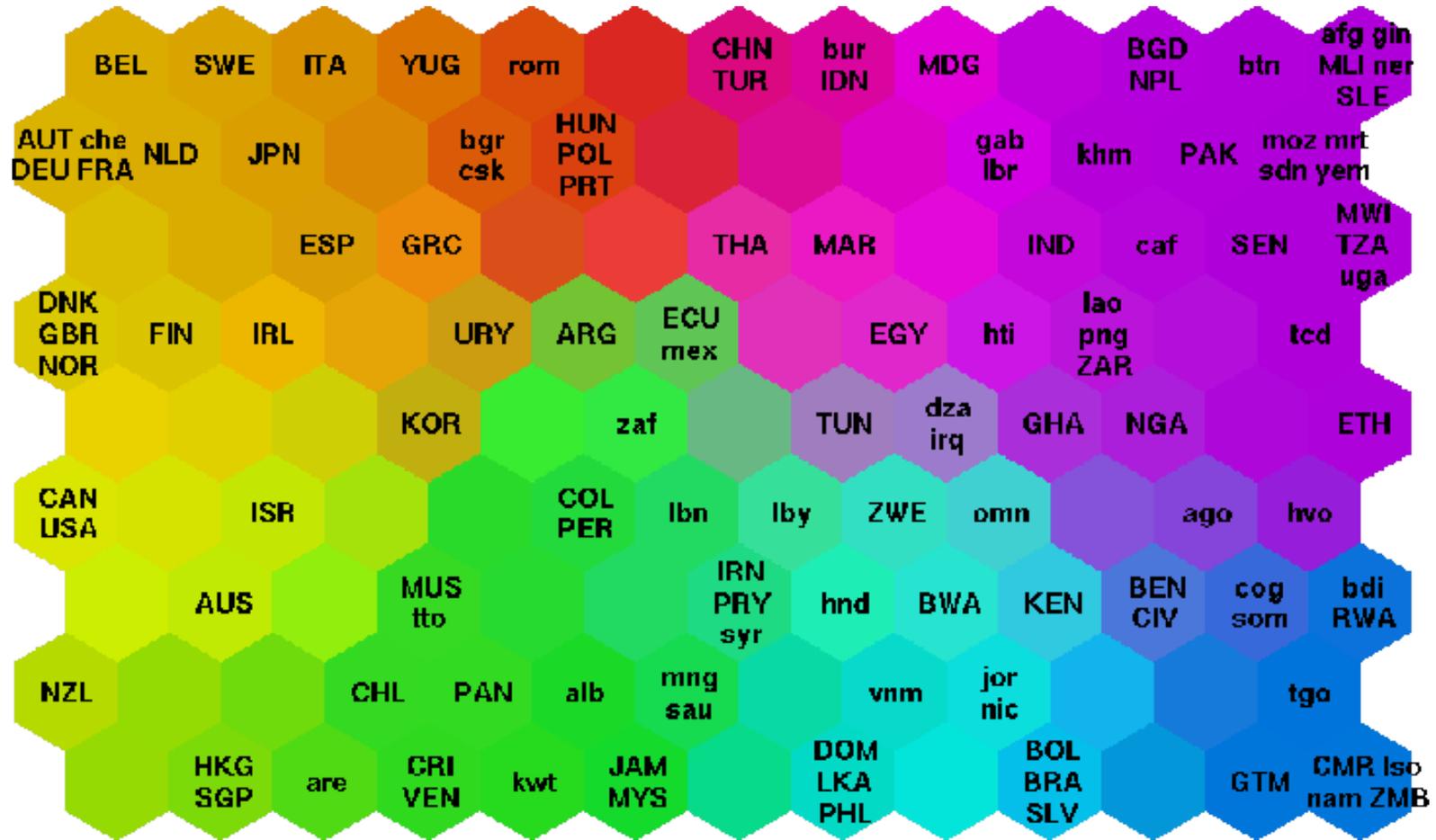
# World Bank Statistics

- ❑ Data: World Bank statistics of countries in 1992.
- ❑ 39 indicators considered e.g., health, nutrition, educational services, etc.
- ❑ The complex joint effect of these factors can be visualized by organizing the countries using the self-organizing map.

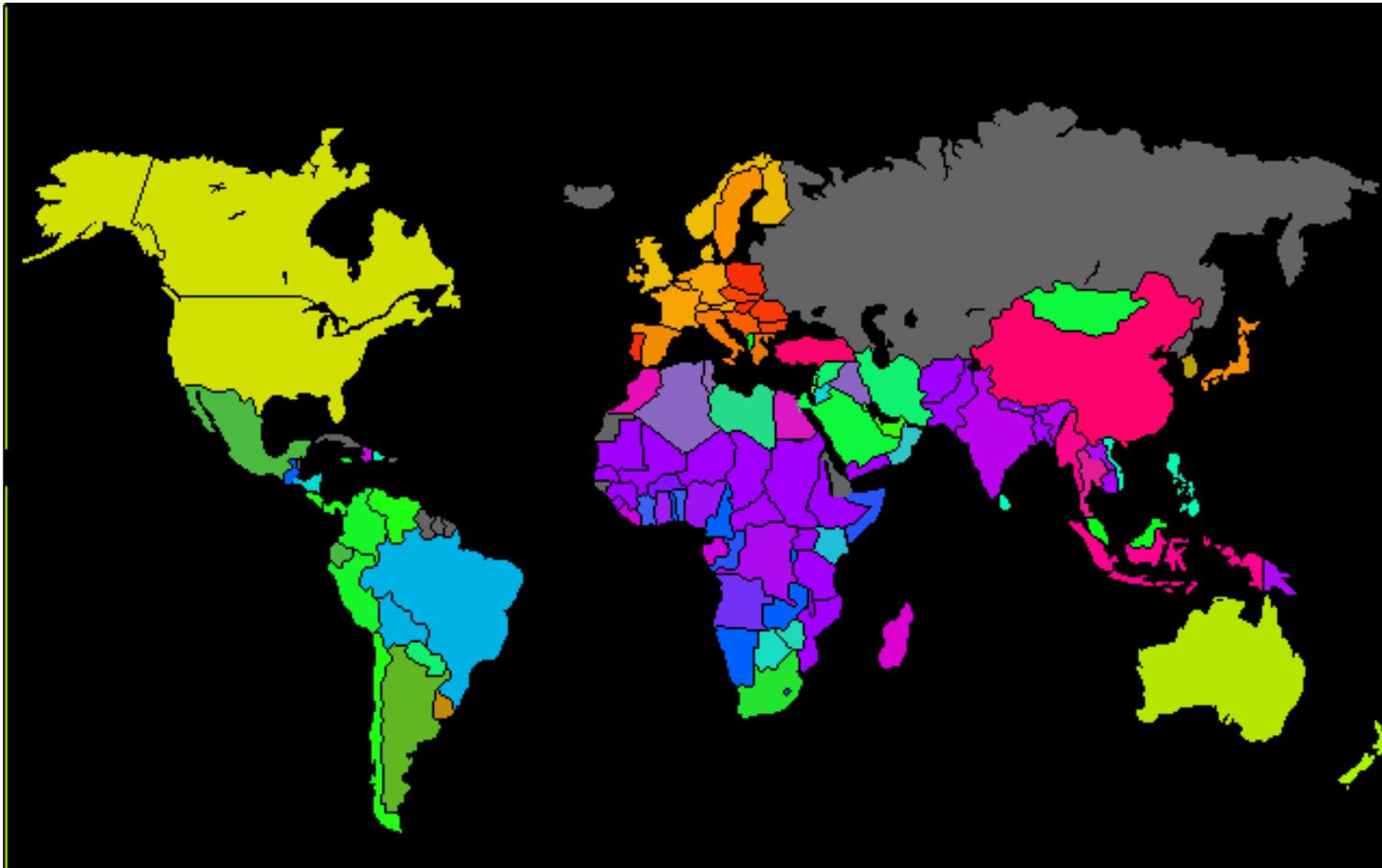
# World Poverty PCA

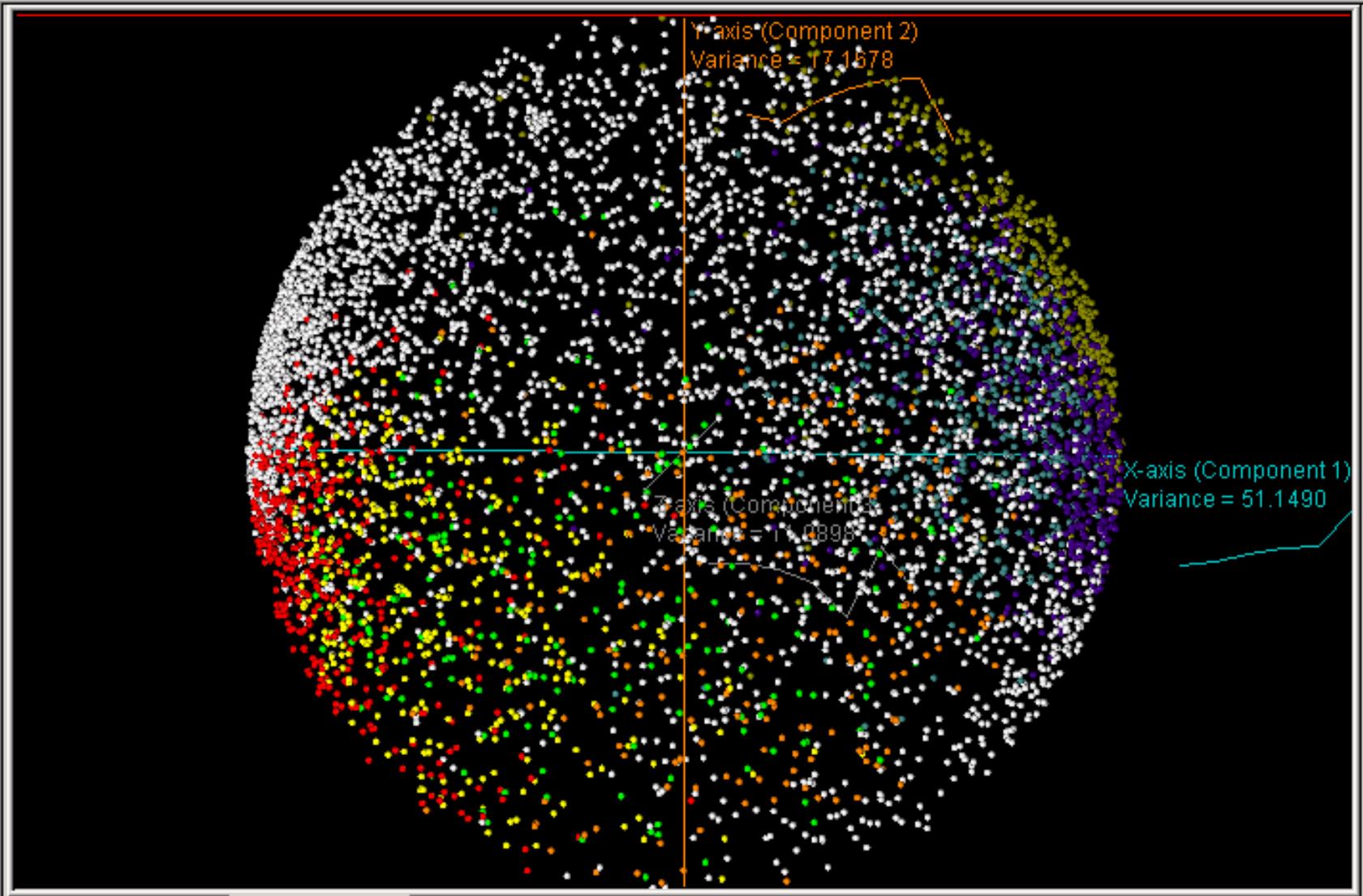


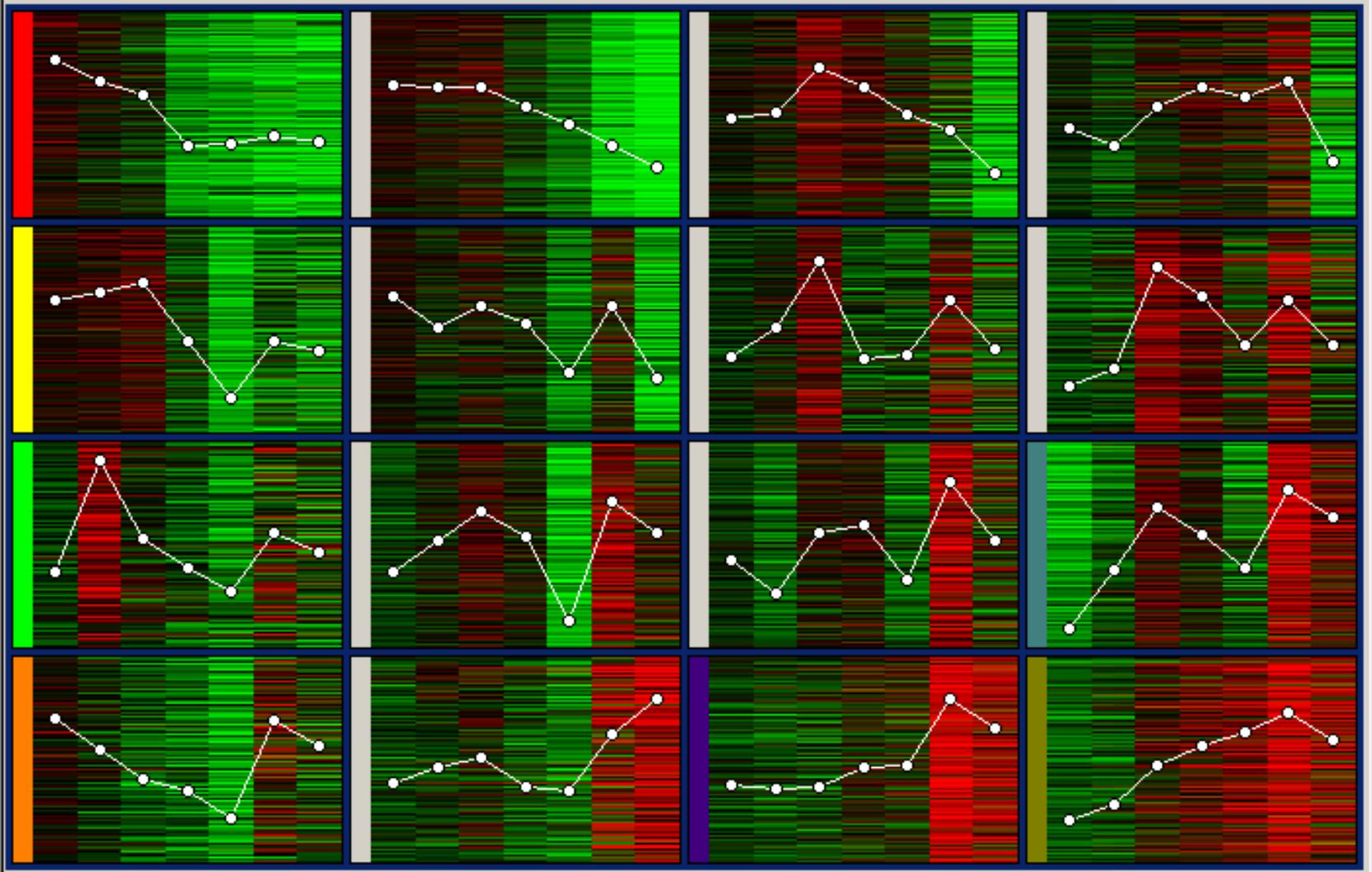
# World Poverty SOM



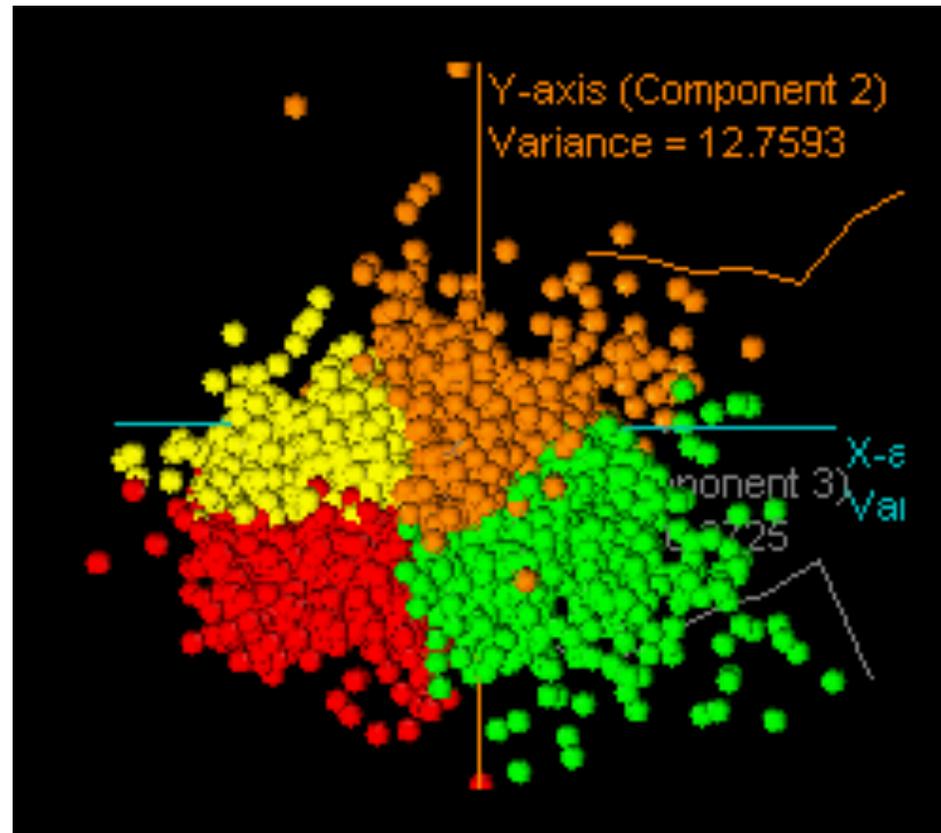
# World Poverty Map



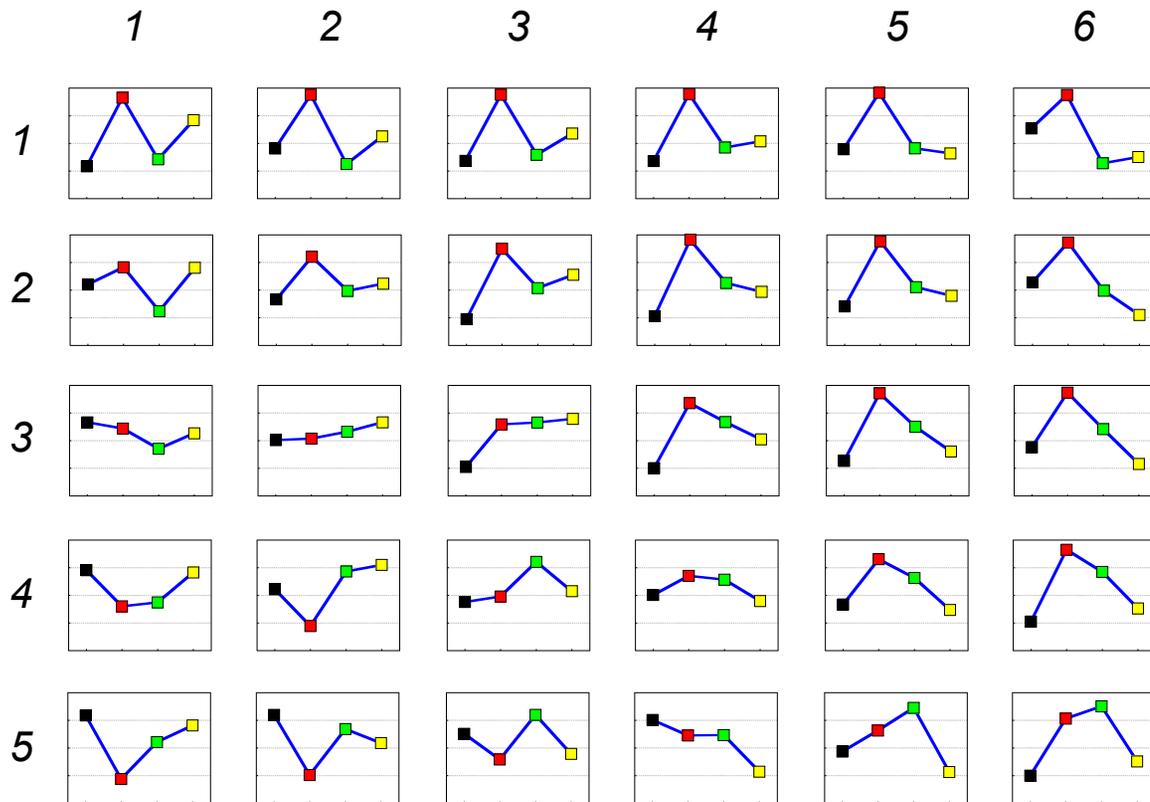




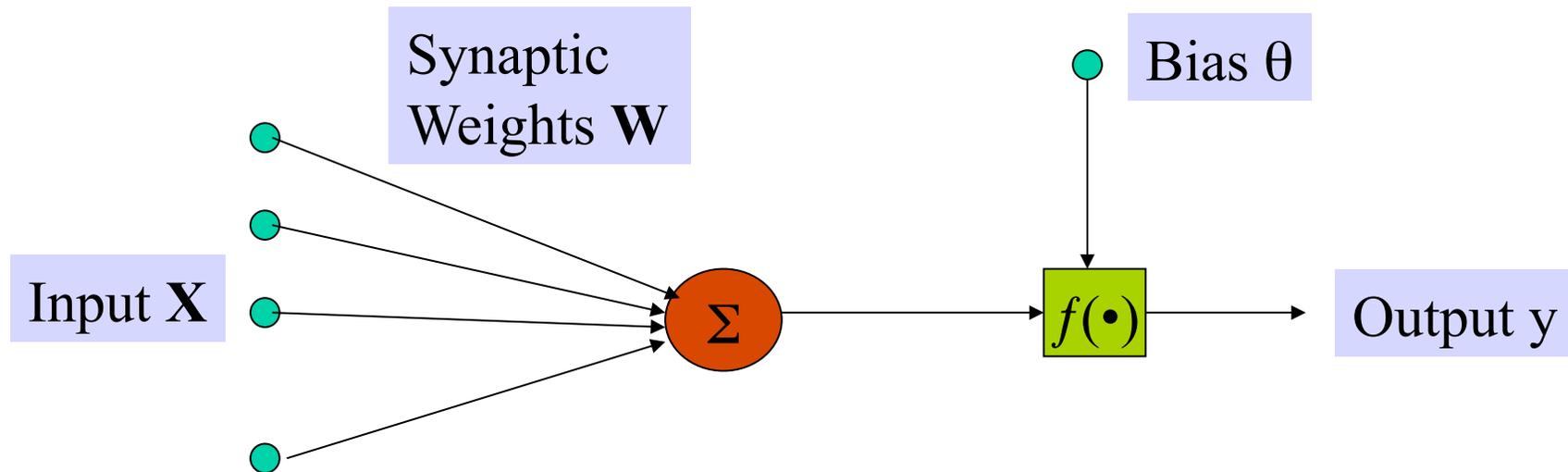
# Viewing SOM Clusters on PCA axes



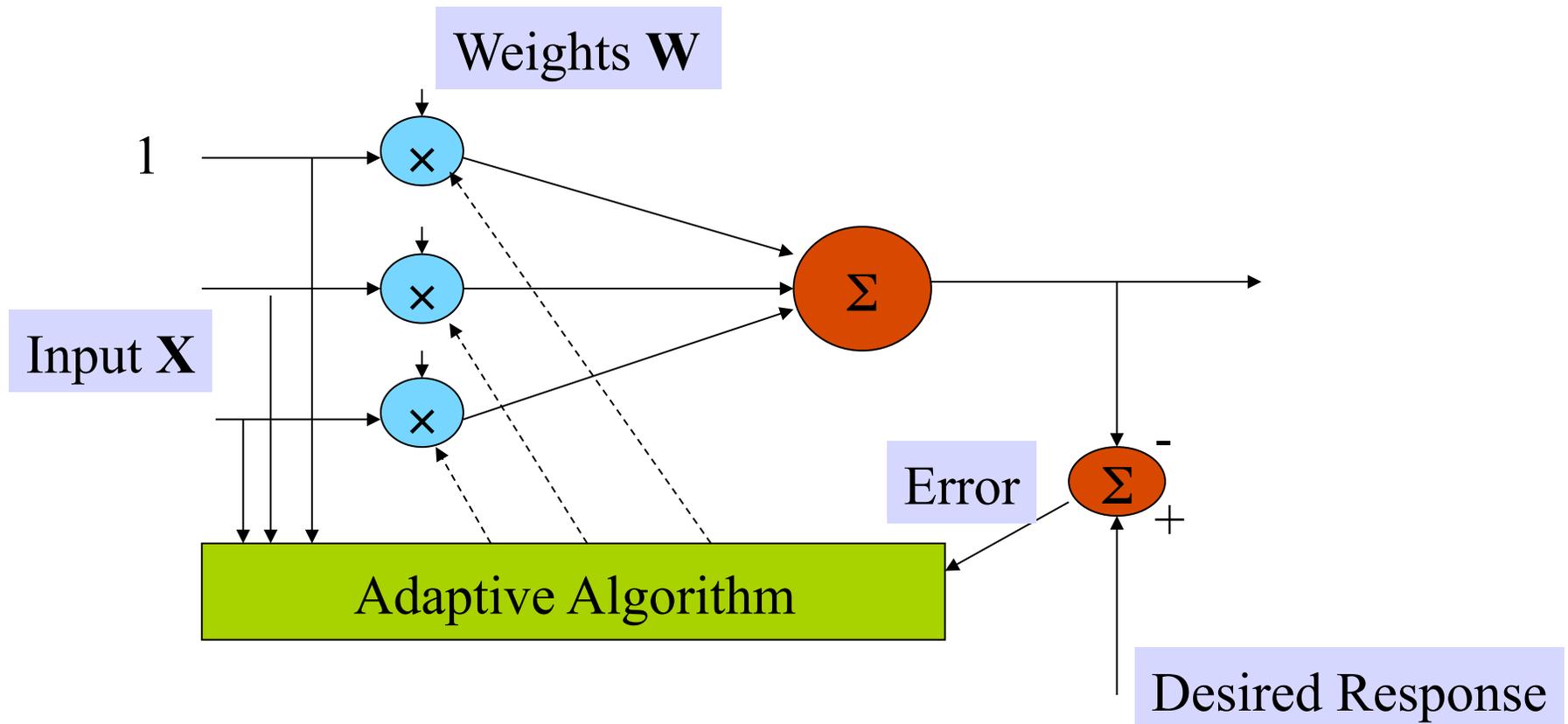
# SOM Example [Xiao-ruì He]



# Neural Networks



# Learning NN



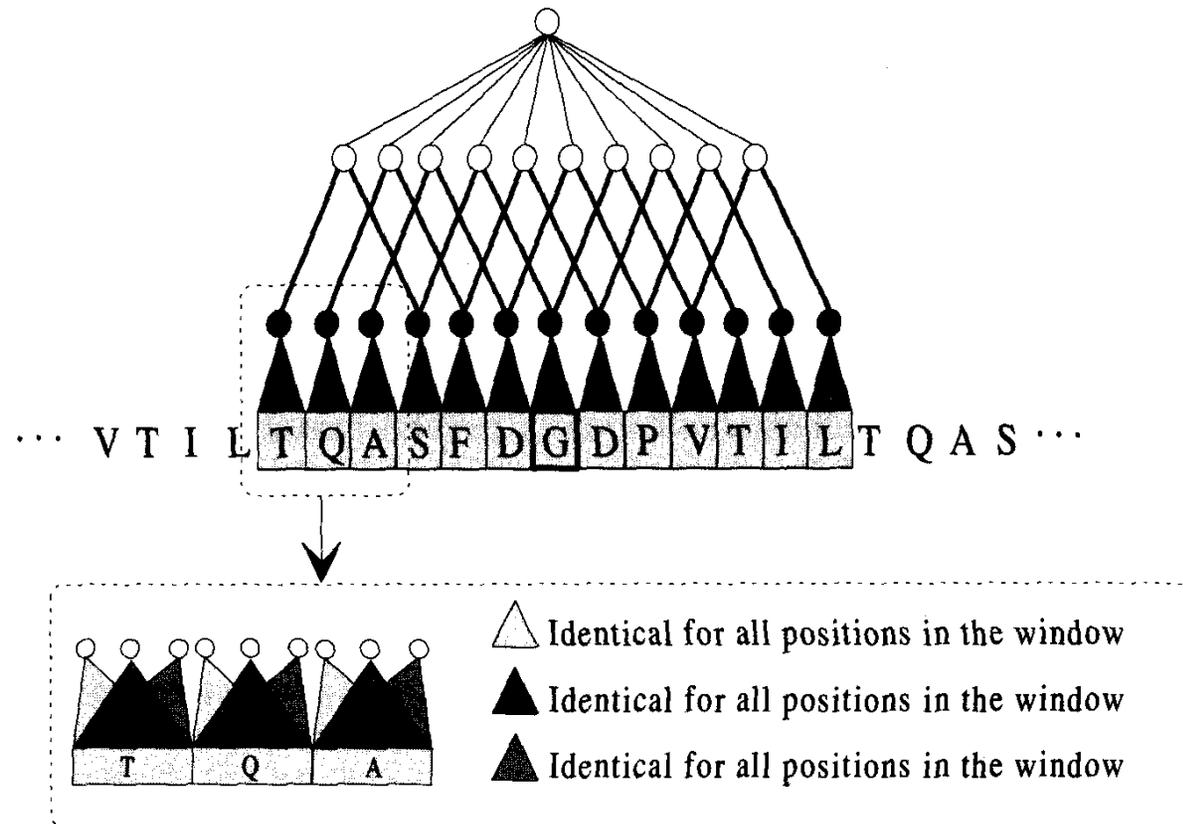
# Types of NNs

- Recurrent NN
- Feed-forward NN
- Layered

# Other issues

- Hidden layers possible
- Different activation functions possible

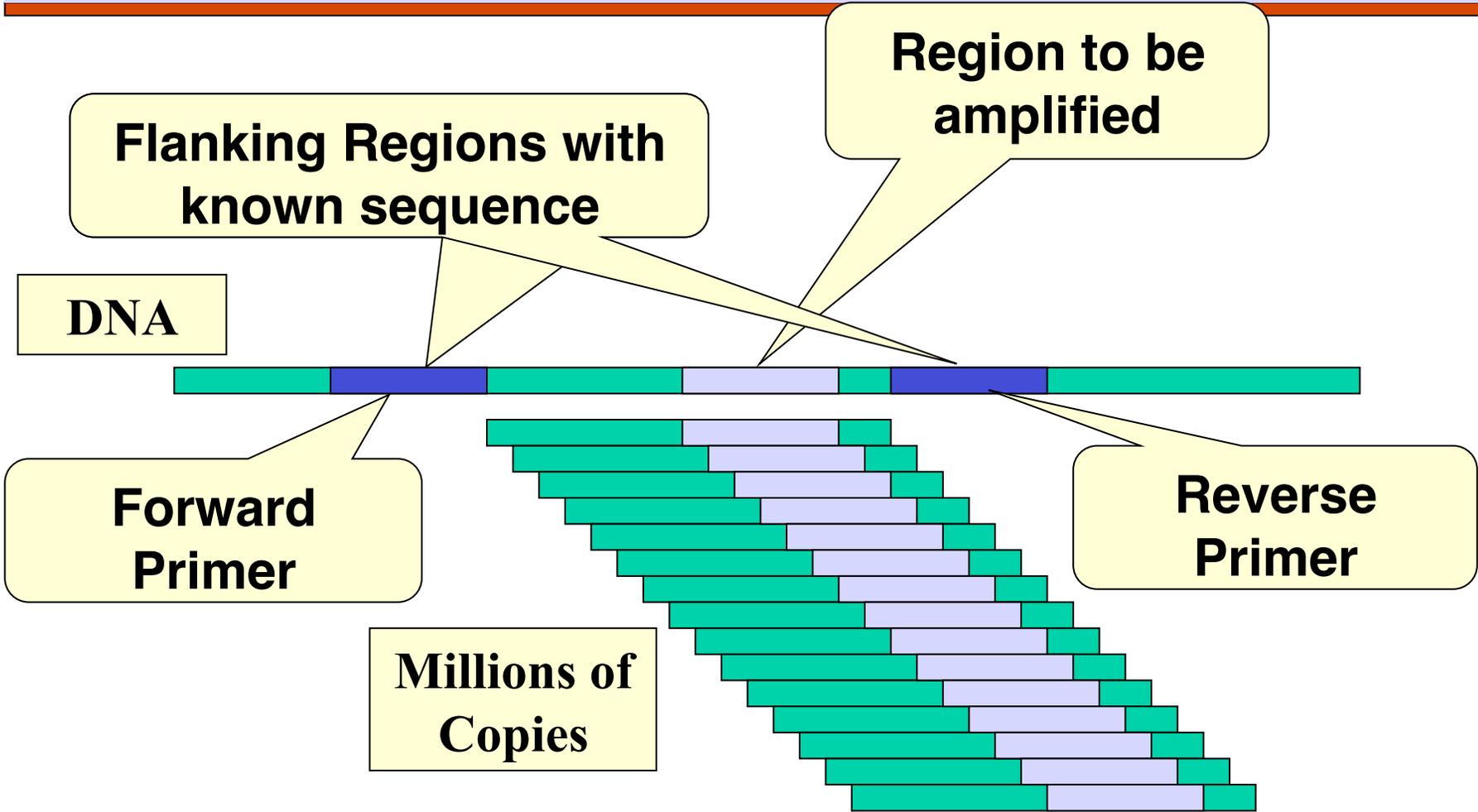
# Application: Secondary Structure Prediction



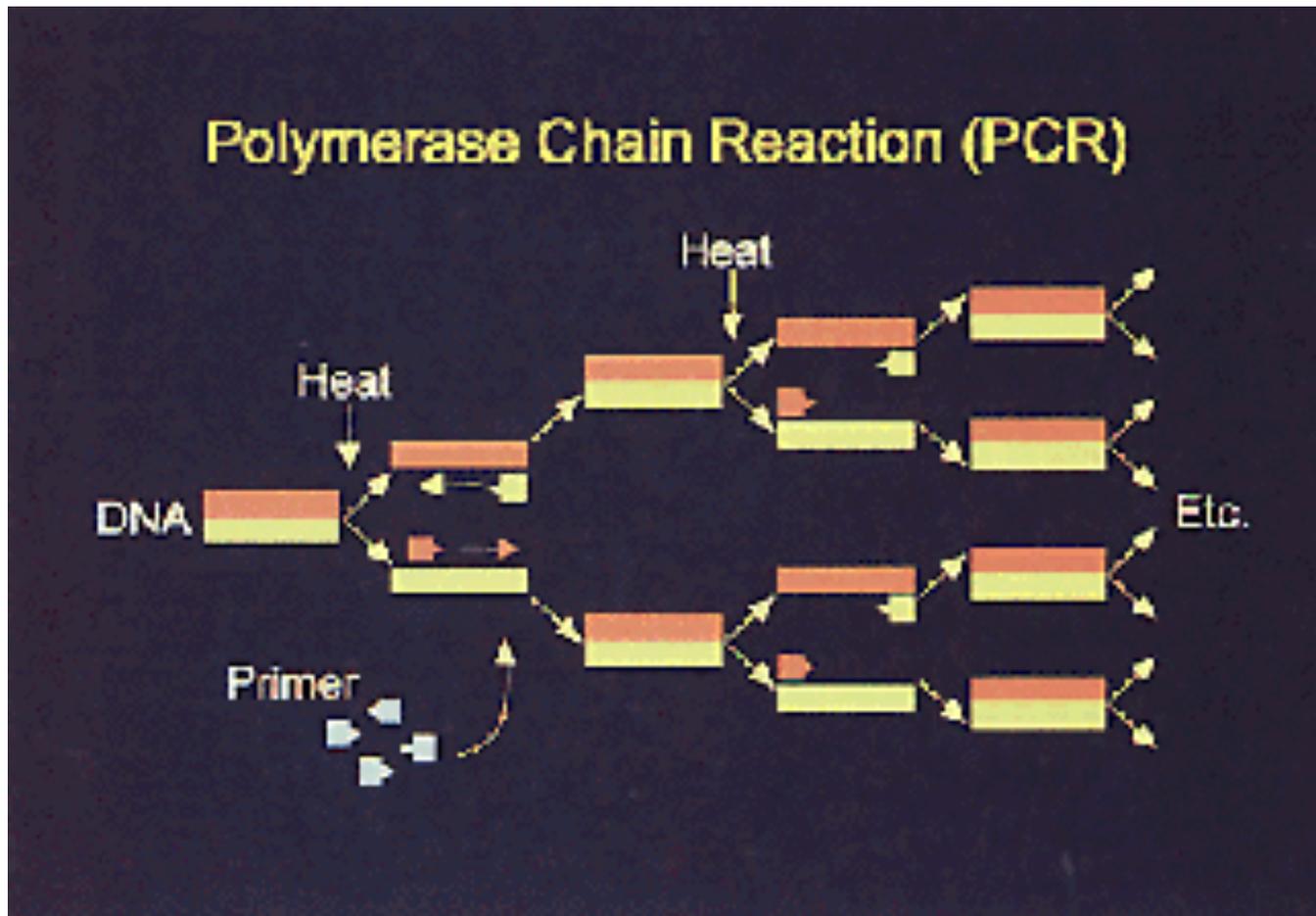
# Polymerase Chain Reaction (PCR)

- ❑ For testing, large amount of DNA is needed
  - Identifying individuals for forensic purposes
    - (0.1 microliter of saliva contains enough epithelial cells)
  - Identifying pathogens (viruses and/or bacteria)
- ❑ PCR is a technique to amplify the number of copies of a specific region of DNA.
- ❑ Useful when exact DNA sequence is unknown
- ❑ Need to know "flanking" sequences
- ❑ Primers designed from "flanking" sequences

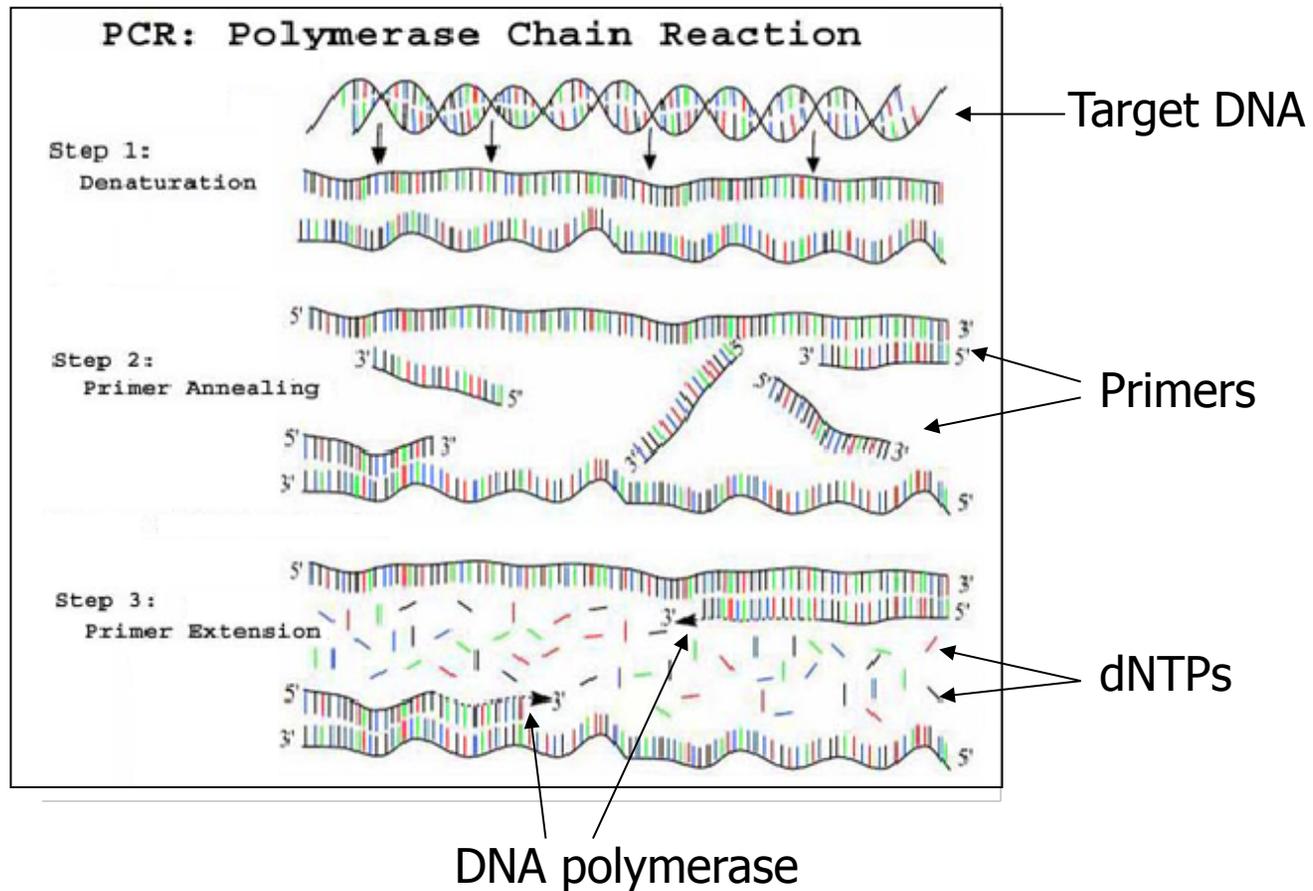
# PCR



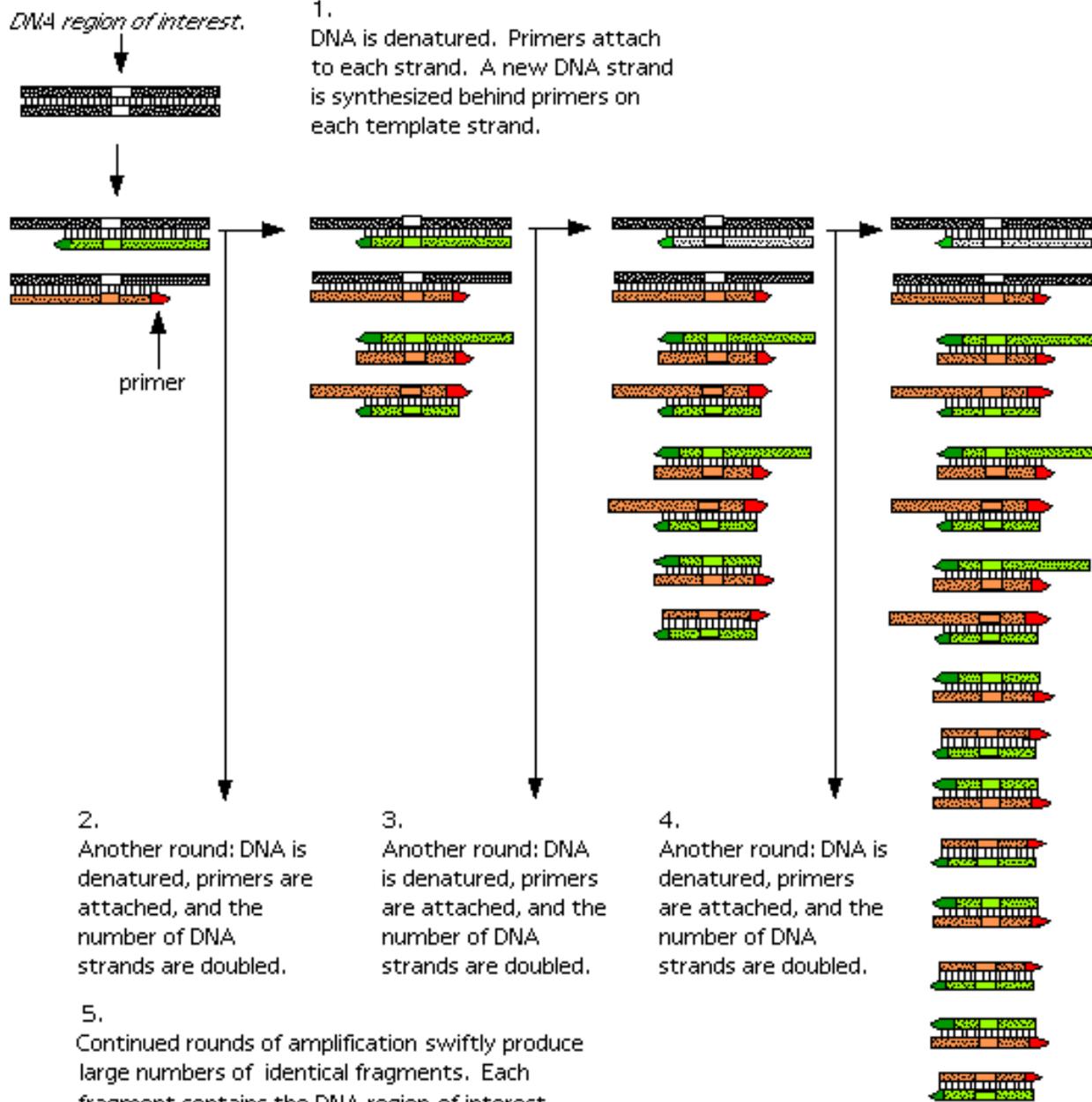
# PCR



# Schematic outline of a typical PCR cycle



# POLYMERASE CHAIN REACTION



1. DNA is denatured. Primers attach to each strand. A new DNA strand is synthesized behind primers on each template strand.

2. Another round: DNA is denatured, primers are attached, and the number of DNA strands are doubled.

3. Another round: DNA is denatured, primers are attached, and the number of DNA strands are doubled.

4. Another round: DNA is denatured, primers are attached, and the number of DNA strands are doubled.

5. Continued rounds of amplification swiftly produce large numbers of identical fragments. Each fragment contains the DNA region of interest.