

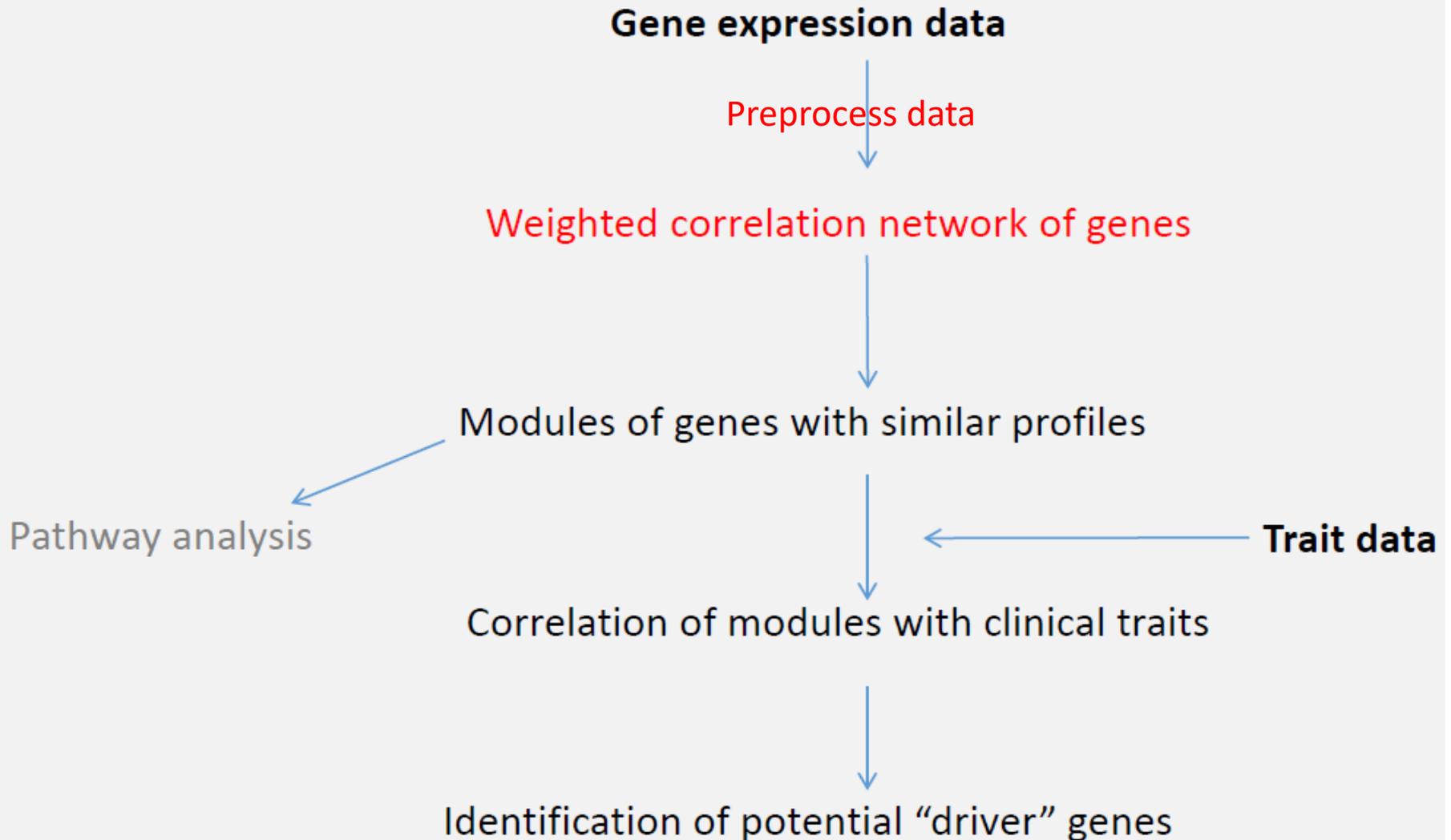
# **Biological Network Analysis and Gene-Gene Interaction Networks - continued**

# **Weighted correlation network analysis (WGCNA)\***

## **\*Citation for WGCNA summary:**

SIB course, Nov. 15-17, 2016, “Introduction to Biological Network Analysis” by Leonore Wigger and with Frédéric Burdet and Mark Ibberson

# Workflow of WGCNA

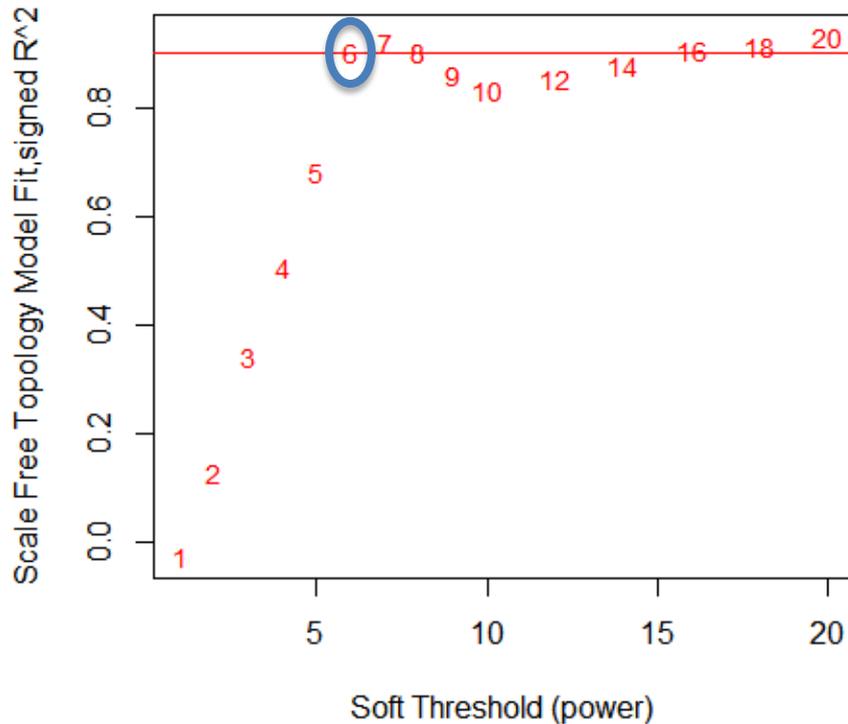


## 4 steps to get **from expression data to adjacency matrix** (explained and detailed last week)

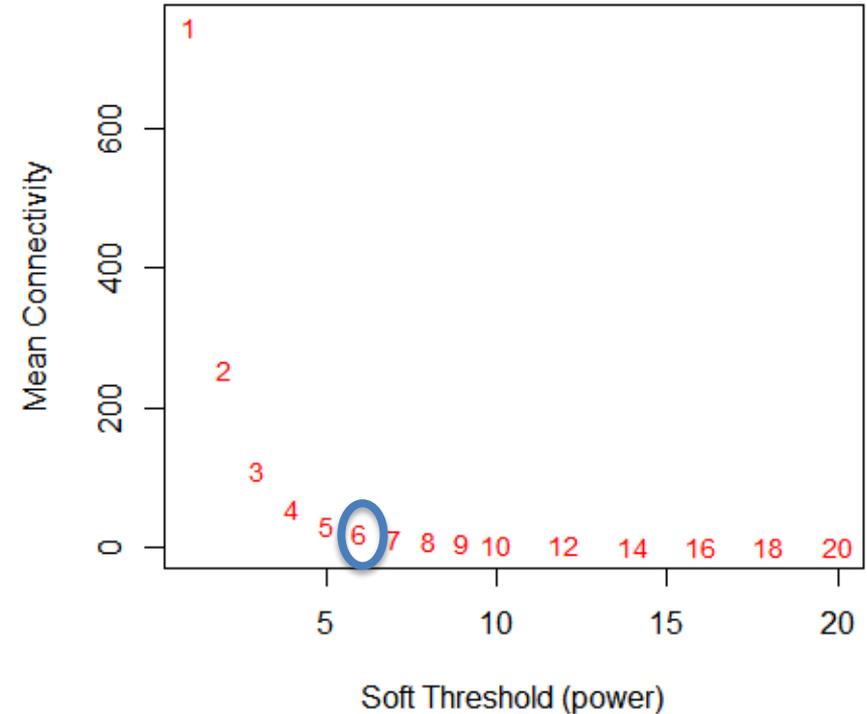
1. Preprocess the gene expression data by cleaning outliers, lowly expressed, and lowly variated genes.
2. Calculate a correlation matrix between gene pairs using e.g. Pearson or Spearman correlation coefficient
3. Pick a soft power value to provide a scale-free topology to the correlation matrix (plot and examine relevant figures for this aim)
4. Calculate the adjacency matrix by applying soft power to the correlation matrix

# Pick a power term: Visual Aid in WGCNA

Scale independence



Mean connectivity



- **Left plot:** Choose power 6. Lowest possible power term where topology approximately fits a scale free network (on or above red horizontal line).
- **Right plot:** mean connectivity drops as power goes up. Must **not** drop too low

# 4 steps to get from network to modules

1. Compute dissimilarity between genes: “topological overlap measure dissimilarity” from the adjacency matrix
2. Perform hierarchical clustering of genes: obtain tree structure
3. Divide clustered genes into modules: cut tree branches
4. **Optional:** Merge very similar modules: use module “eigengenes”

# Step 1: Compute dissimilarity between genes

- **Why we use Topological Overlap Measure (TOM)?**
  - TOM is a pairwise similarity measure between network nodes (genes)
  - TOM(i,j) is **high if genes i,j have many shared neighbors** because overlap of their network neighbors is large
- **So, a high TOM(i,j) implies that genes have similar expression patterns**
- **How to calculate TOM similarity between two nodes:**
  - 1. Count number of shared neighbors: “agreement” of the set of neighboring nodes
  - 2. Normalize to [0,1]

TOM(i,j) = 0 means: no overlap of network neighbors

TOM(i,j) = 1 means: identical set of network neighbors

- $a_{ij}$  is the adjacency score between genes i and j taken from the adjacency matrix

$$TOM_{ij} = \frac{\sum_u a_{iu} a_{uj} + a_{ij}}{\min(k_i, k_j) + 1 - a_{ij}}$$
$$DistTOM_{ij} = 1 - TOM_{ij}$$

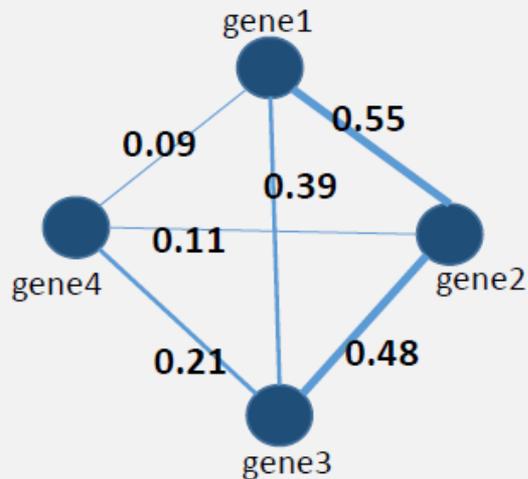
- ***NOTE that: Generalized to the case of weighted networks in Zhang and Horvath (2005), first WGCNA paper***
  - All nodes are neighbors; counting them is not informative.
  - Compute agreement of the set of neighboring nodes based on edge strengths.

- But, we need a **dissimilarity measure** for clustering!!
- TOM as a **similarity measure** can be transformed into a **dissimilarity measure**: **distTOM= 1-TOM.**

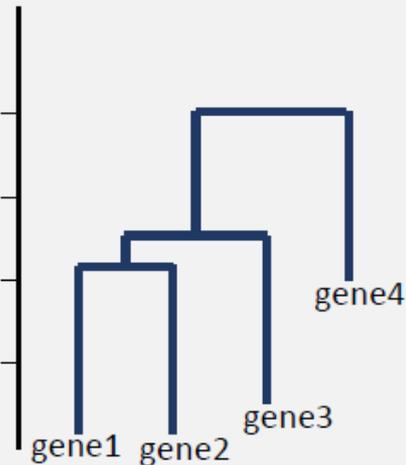
## Step 2: Perform hierarchical clustering of genes

- **Compute gene dendrogram:**

Weighted correlation network  
from gene expression data



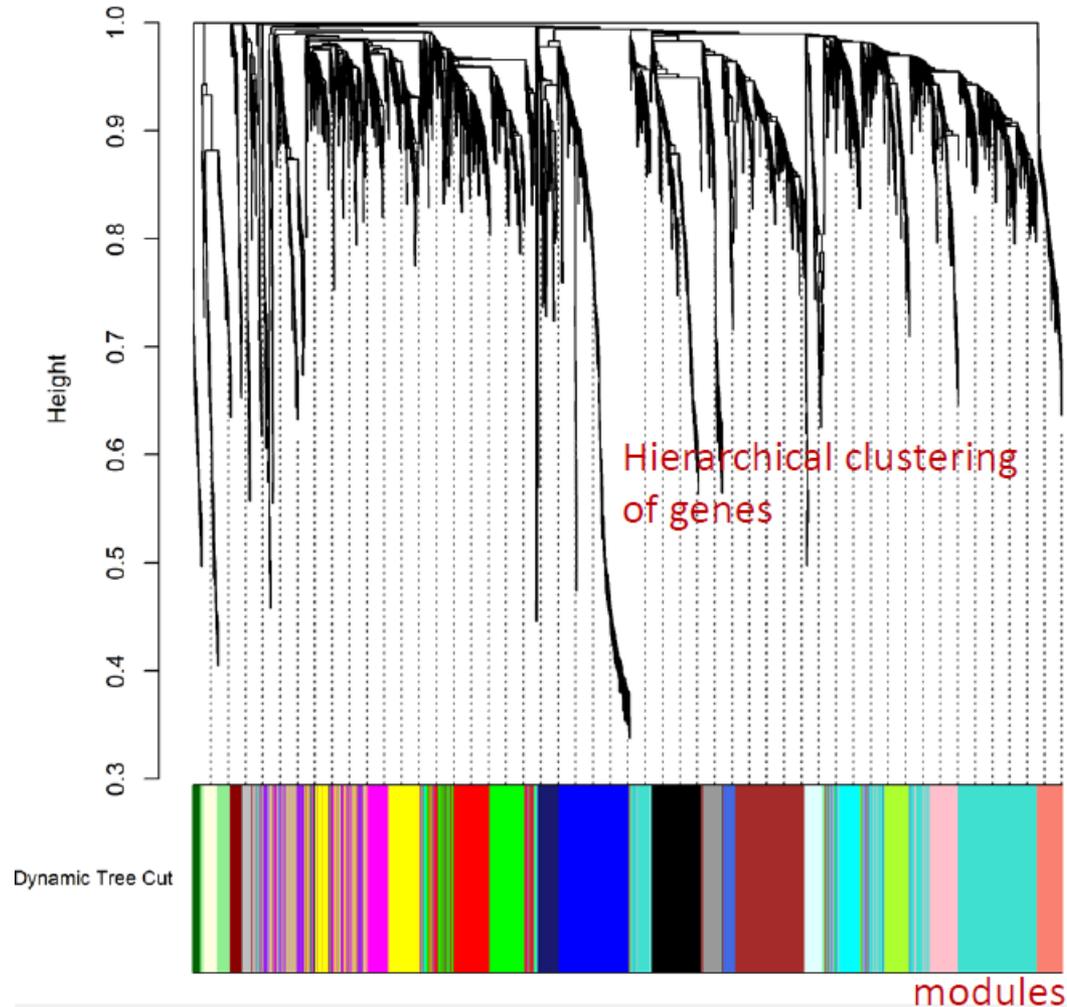
Clustering dendrogram



*(dis)similarity between genes:  
Topological Overlap Measure TOM*

# Step 3: Divide clustered genes into modules

- Gene dendrogram and detected modules

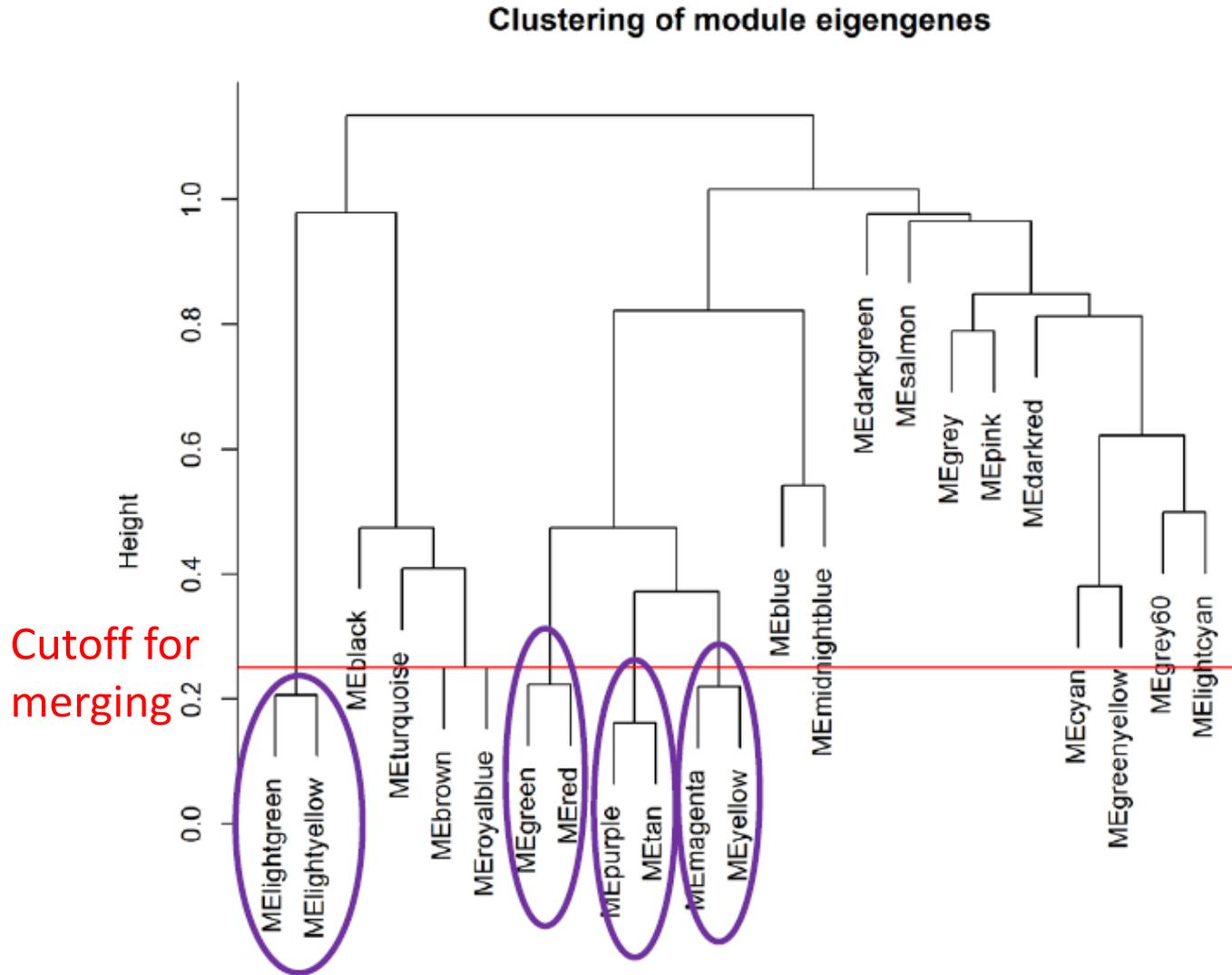


Dynamic tree cut algorithm groups genes into modules:  
*corFnc="pearson"; power=6; min.modulesize=30*

## Step 4 (Optional): Merge very similar modules

- **A module eigengene** is a 1-dimensional data vector, summarizing the expression data of the genes that form a module
- **How it is computed:** the 1st principal component of the expression data
- **What it is used for:** to represent the module in mathematical operations:
  - modules can be correlated with one another
  - modules could be **clustered together** (we can combine them)
  - modules can be correlated with external traits

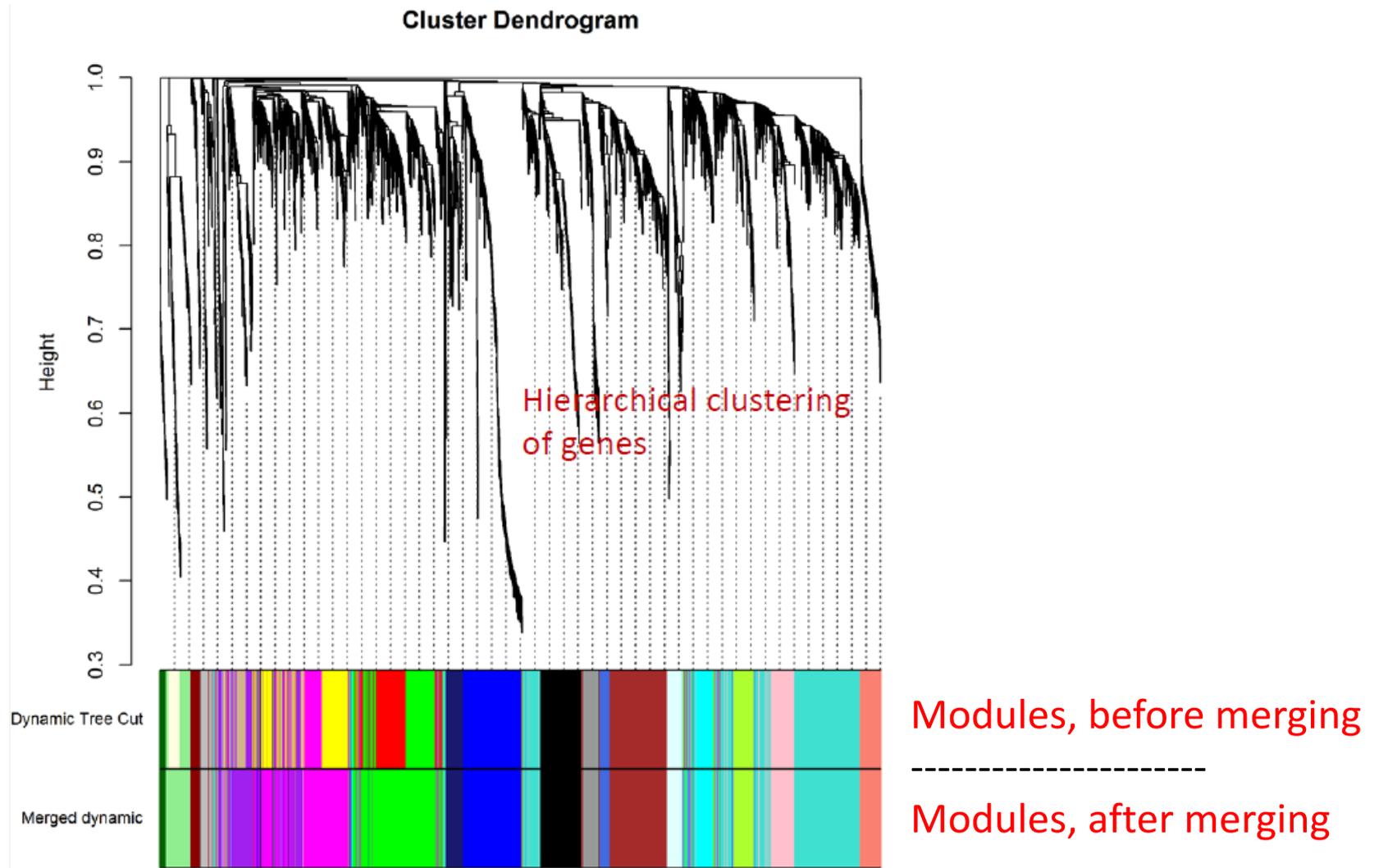
# Clustering of module eigengenes



***Dissimilarity measure:  $1 - \text{cor}(\text{Meigengenes})$***

**Merge modules whose dissimilarity is below the merging cutoff**

# Gene dendrogram and detected modules, before and after merging



*corFnc="pearson"; power=6; min.modulesize=30*

# Module Detection: Decisions to make

## How to choose optimal parameters?

### Dynamic Tree cut procedure

- **Minimal module size:** typically 20 or 30
- **CutHeight:** try different heights such as 0.95 or 0.9995, etc.

### Module Merging procedure

- **Cutoff for module eigengene dendrogram:** typically between 0.15 and 0.25
  - *check if clusters look ok on dendrogram*
- **Merge once or several times?** Usually once, but merge step can be repeated:
  - *if some modules are very similar*
  - *if we want larger modules*

# Module preservation analysis in WGCNA

**Is my network module preserved and reproducible?\***

\* Langfelder et al PloS Comp Biol. 7(1): e1001057.

# Network-based module preservation statistics

- Input: module assignment in reference data.
- Adjacency matrices in **reference**  $A^{\text{ref}}$  and **test** data  $A^{\text{test}}$
- Network preservation statistics assess preservation of
  - 1. network density: Does the module remain densely connected in the test network?
  - 2. connectivity: Is hub gene status preserved between reference and test networks?
  - 3. separability of modules: Does the module remain distinct in the test data?

# Several connectivity preservation statistics

*For general networks, i.e. input adjacency matrices*

- $cor.kIM = cor(kIM^{ref}, kIM^{test})$ 
  - *correlation of intramodular connectivity across module nodes*
- $cor.ADJ = cor(A^{ref}, A^{test})$ 
  - *correlation of adjacency across module nodes*

For correlation networks, i.e. input sets are variable measurements

- $cor.Cor = cor(cor^{ref}, cor^{test})$
- $cor.kME = cor(kME^{ref}, kME^{test})$

One can derive relationships among these statistics in case of weighted correlation network

# Choosing thresholds for preservation statistics based on permutation test

- For correlation networks, we study **4 density** and **3 connectivity** preservation statistics that take on values  $\leq 1$
- Challenge: Thresholds could depend on many factors (number of genes, number of samples, biology, expression platform, etc.)
- Solution: Permutation test. Repeatedly permute the gene labels in the test network to estimate the mean and standard deviation under the null hypothesis of no preservation.
- Next we calculate a **Z statistic**:  $Z = (\text{observed} - \text{mean}) / \text{sd}$
- We have had **4 density** and **3 connectivity** preservation statistics:

$$Z_{density} = \text{median}(Z_{meanCor}, Z_{meanAdj}, Z_{propVarExpl}, Z_{meanKME}).$$

$$Z_{connectivity} = \text{median}(Z_{cor.kIM}, Z_{cor.kME}, Z_{cor.cor}).$$

# Permutation test for estimating Z scores

- For each preservation measure we report the observed value and the permutation Z score to measure significance ( $Z = (\text{observed} - \text{mean}) / \text{sd}$ ).
- Each Z score provides answer to “Is the module significantly better than a random sample of genes?”
- Summarize the individual Z scores into a composite measure called **Z.summary**

$$Z_{summary} = \frac{Z_{density} + Z_{connectivity}}{2}.$$

- **Z.summary < 2 indicates no preservation,**
- **2 < Z.summary < 10 weak to moderate evidence of preservation,**
- **Z.summary > 10 strong evidence**

# Summary preservation

- Standard cross-tabulation based statistics are intuitive
  - Disadvantages: i) only applicable for modules defined via a module detection procedure, ii) ill suited for ruling out module preservation
- Network based preservation statistics measure different aspects of module preservation
  - Density-, connectivity-, separability preservation
- Two types of composite statistics: **Zsummary** and **medianRank**.
- Composite statistic **Zsummary** based on a permutation test
  - Advantages: thresholds can be defined, R function also calculates corresponding permutation test p-values
  - Example:  $Zsummary < 2$  indicates that the module is *\*not\** preserved
  - Disadvantages: i) Zsummary is computationally intensive since it is based on a permutation test, ii) often depends on module size
- Composite statistic **medianRank**
  - Advantages: i) fast computation (no need for permutations), ii) no dependence on module size.
  - Disadvantage: only applicable for ranking modules (i.e. relative preservation)

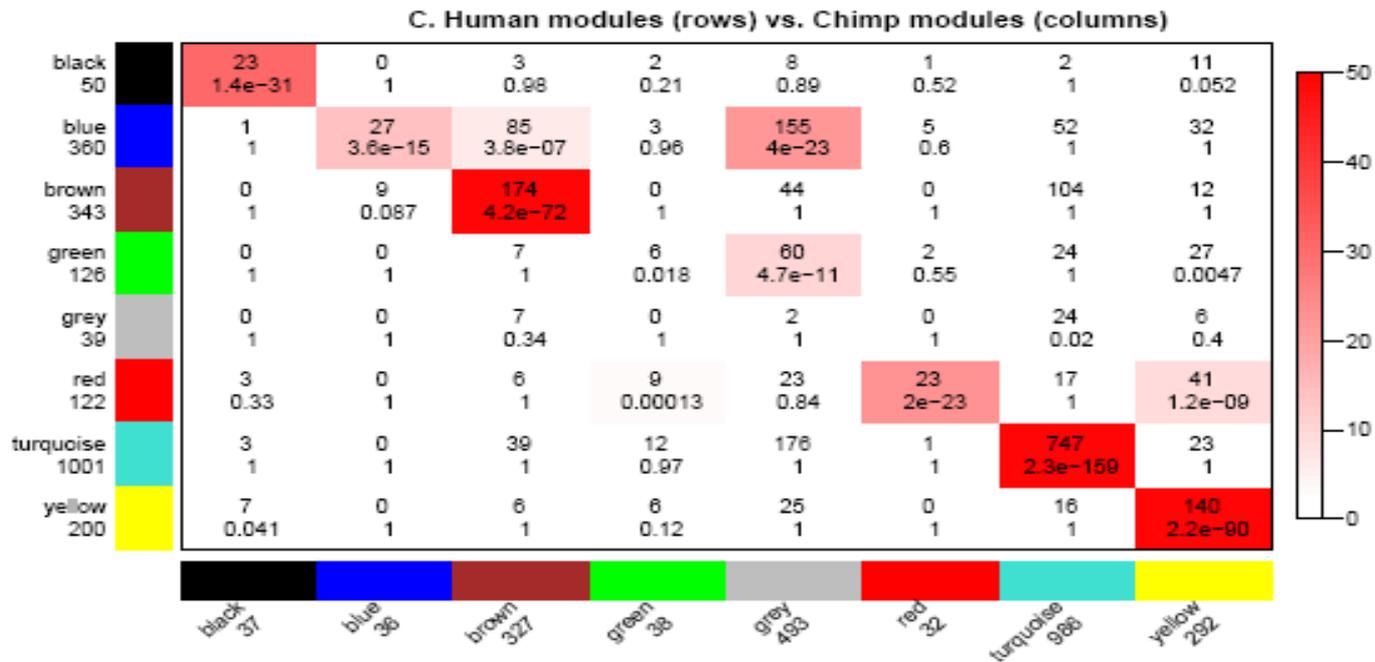
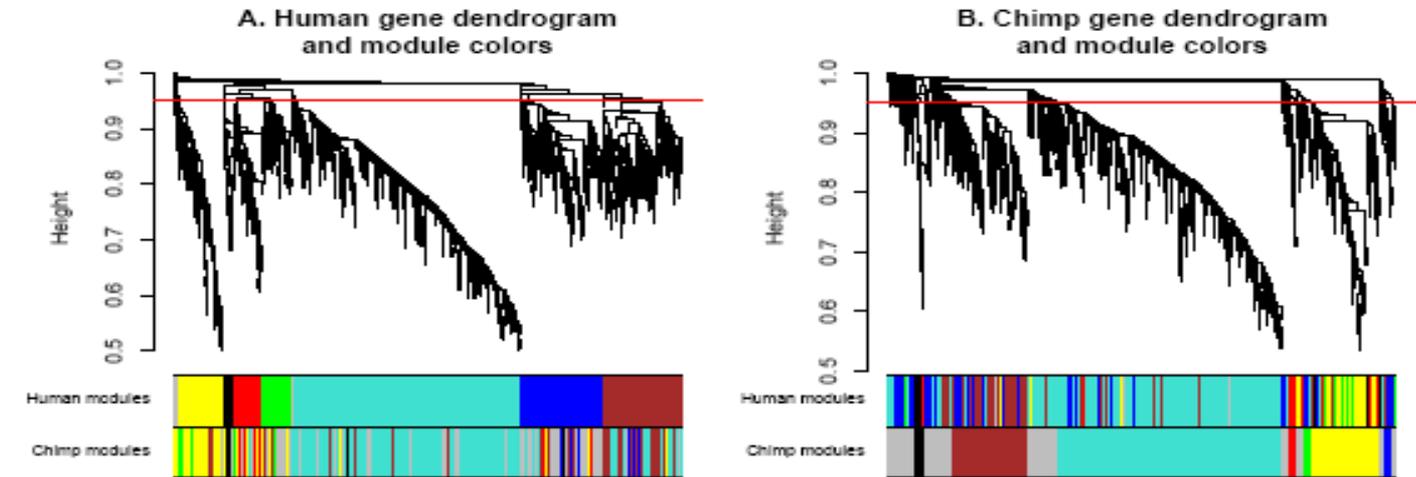
## **Application:**

**Studying the preservation of human brain co-expression modules in chimpanzee brain expression data.**

Modules defined as clusters  
(branches of a cluster tree)

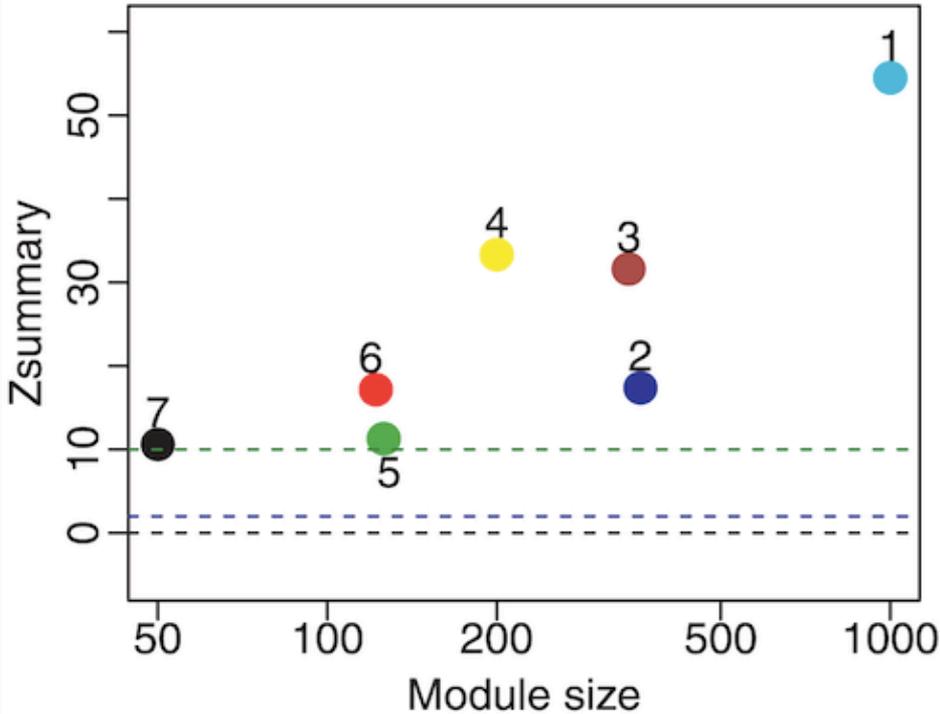
Data from Oldam et al 2006

# Preservation of modules between human and chimpanzee brain networks

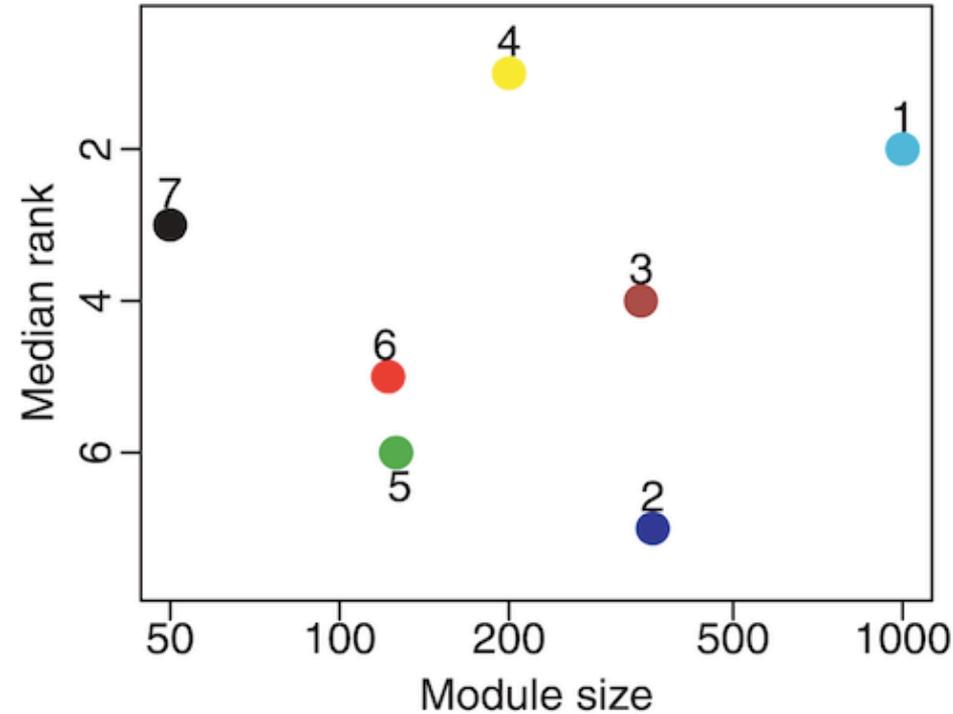


# 2 composite preservation statistics

A. Zsummary



B. Median rank



- Zsummary is above the threshold of 10 (green dashed line), i.e. all modules are preserved.
- Zsummary often shows a dependence on module size which may or may not be attractive
- In contrast, the median rank statistic is not dependent on module size.
- It indicates that the yellow module is the most preserved one

# Implementation and R software tutorials, WGCNA R library

- General information on weighted correlation networks
- Google search
  - “WGCNA”
  - “weighted gene co-expression network”
- R function `modulePreservation` is part of WGCNA package
- Tutorials: preservation between human and chimp brains

[www.genetics.ucla.edu/labs/horvath/CoexpressionNetwork/ModulePreservation](http://www.genetics.ucla.edu/labs/horvath/CoexpressionNetwork/ModulePreservation)

# Pseudocode for modulePreservation()

- `setLabels = c("Control", "Disease");`
- `multiExpr = list(Control = list(data = dataExp_Control), Disease = list(data = dataExp_Disease));`
- `multiLabel = list(Control = moduleLabelsControl, Disease=moduleLabelsDisease);`
- `mp = modulePreservation(multiExpr, multiLabel, referenceNetworks = c(1:2), nPermutations = 100, randomSeed = 1, quickCor = 0, verbose =3);`
- `save(multiExpr, multiLabel, mp, file =“modulepreservationLabel_Alldata.RData” )`
- `ref = 1; test = 2`
- `Zsummary1=mp$preservation$Z[[ref]][[test]][, 2]`
- `names(Zsummary1)=rownames(mp$preservation$Z[[ref]][[test]])`
- `low.preserved1=Zsummary1[which(Zsummary1<2)]`
- `ref=2; test=1`
- `Zsummary2=mp$preservation$Z[[ref]][[test]][, 2]`
- `names(Zsummary2)=rownames(mp$preservation$Z[[ref]][[test]])`
- `low.preserved2=Zsummary2[which(Zsummary2<2)]`

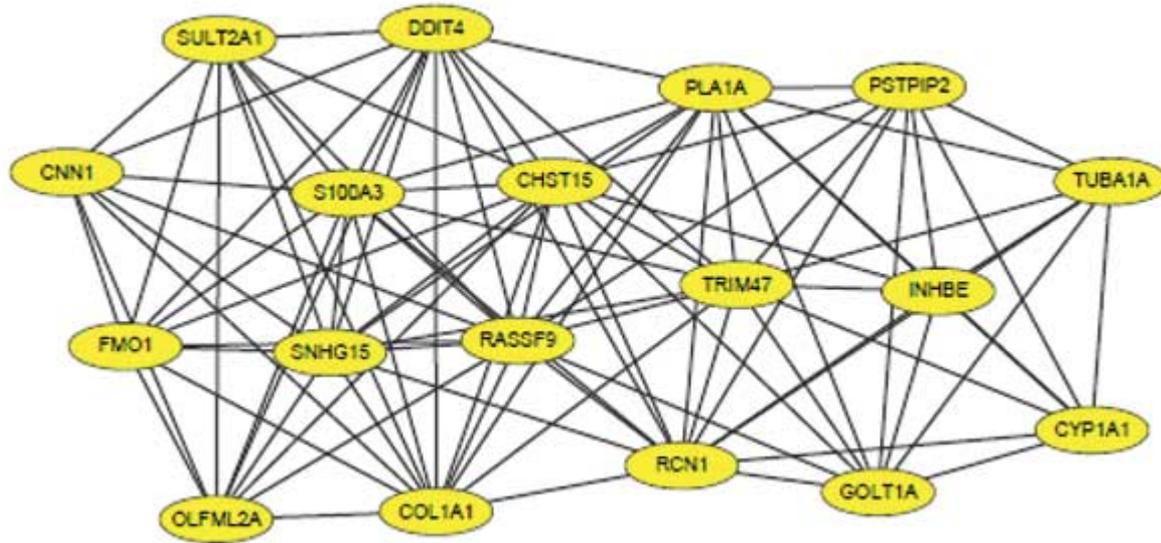
# Introduction to Bayesian Networks

# Content

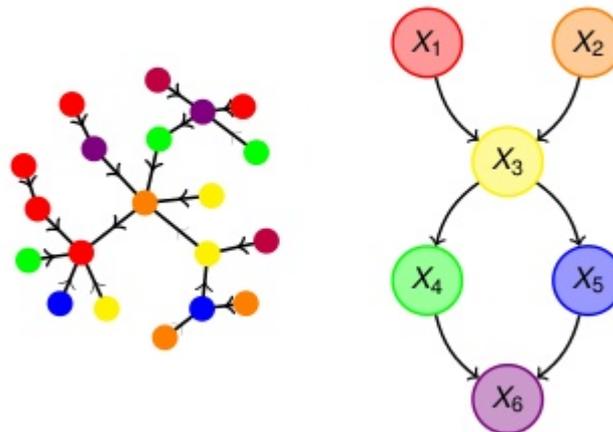
- Definition of Bayesian networks (BN)
- An example, Rimbanet...
- Mandatory and optional input files
- Comparison of BN tools

# Gene-gene interaction networks

## Co-expression networks

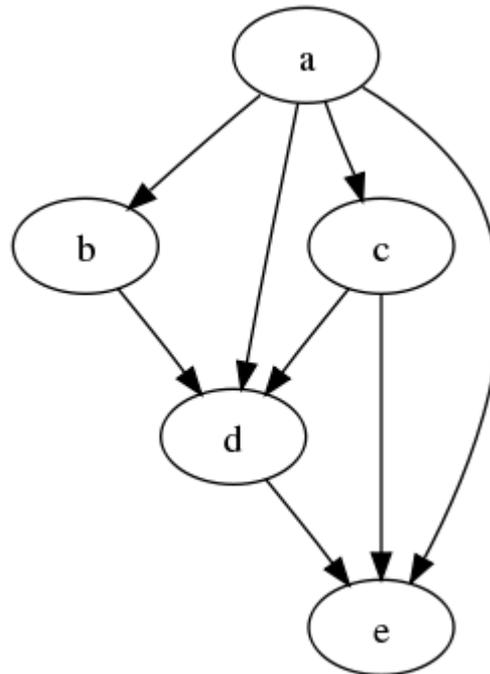


## Bayesian networks



# What is a DAG

- A directed acyclic graph (DAG) is a finite directed graph with no directed cycles.
- There is no way to start at any vertex  $v$  and follow a consistently-directed sequence of edges that eventually loops back to  $v$  again.



# Bayes' theorem

- Bayes' theorem describes the probability of an event, based on prior knowledge of conditions that might be related to the event.
- Bayes' theorem is stated as:

$$P(A, B) = P(B | A)P(A) = P(A | B)P(B)$$
$$\Rightarrow P(A | B) = \frac{P(B | A)P(A)}{P(B)} = \frac{P(B | A)P(A)}{\sum_i P(B | A_i)P(A_i)}$$

# Bayes' theorem - an example

- Suppose a test to detect a disease in human is 99% sensitive and 95% specific, means:

For the patients, the test gives a positive result with 99% (TP) and gives a negative result with 1% (FN):  $P(+ | C)=0.99$  and  $P(- | C)=0.01$

For the healthy samples, test gives a negative result with 95% (TN) and a positive result with 5% (FP):  $P(- | H)=0.95$  and  $P(+ | H)=0.05$

- **Q: If a randomly selected individual's test result is positive, what is the probability that he/she carries the disease? Hence,  $P(C | +) = ?$**

$$\begin{aligned} P(C|+) &= \frac{P(+|C)P(C)}{P(+)} = \frac{P(+|C)P(C)}{\sum_i P(+|Option_i)P(Option_i)} \\ &= \frac{P(+|C)P(C)}{P(+|C)P(C) + P(+|H)P(H)} = \frac{0.99 \times 0.02}{0.99 \times 0.02 + 0.05 \times 0.98} \approx 29\% \end{aligned}$$

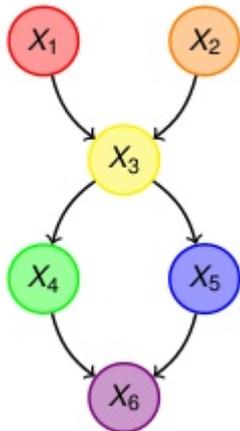
# Bayesian networks (BNs)

- BNs are directed acyclic graphs (DAGs) in which the edges of the graph are defined by conditional probabilities that characterize the distribution of the states of each node given the state of its parents.

$P(M|D) \approx P(D|M) \times P(M)$  , where M: network model, D: observation data

- We can define a partitioned joint probability distribution over all nodes:

$P(X) = \prod_i P(X^i | \text{Pa}(X^i))$  where  $\text{Pa}(X^i)$  is parent set of the node  $X^i$ .



$$P(X_1, \dots, X_6) = P(X_1)P(X_2)P(X_3|X_1, X_2) \\ P(X_4|X_3)P(X_5|X_3)P(X_6|X_4, X_5)$$

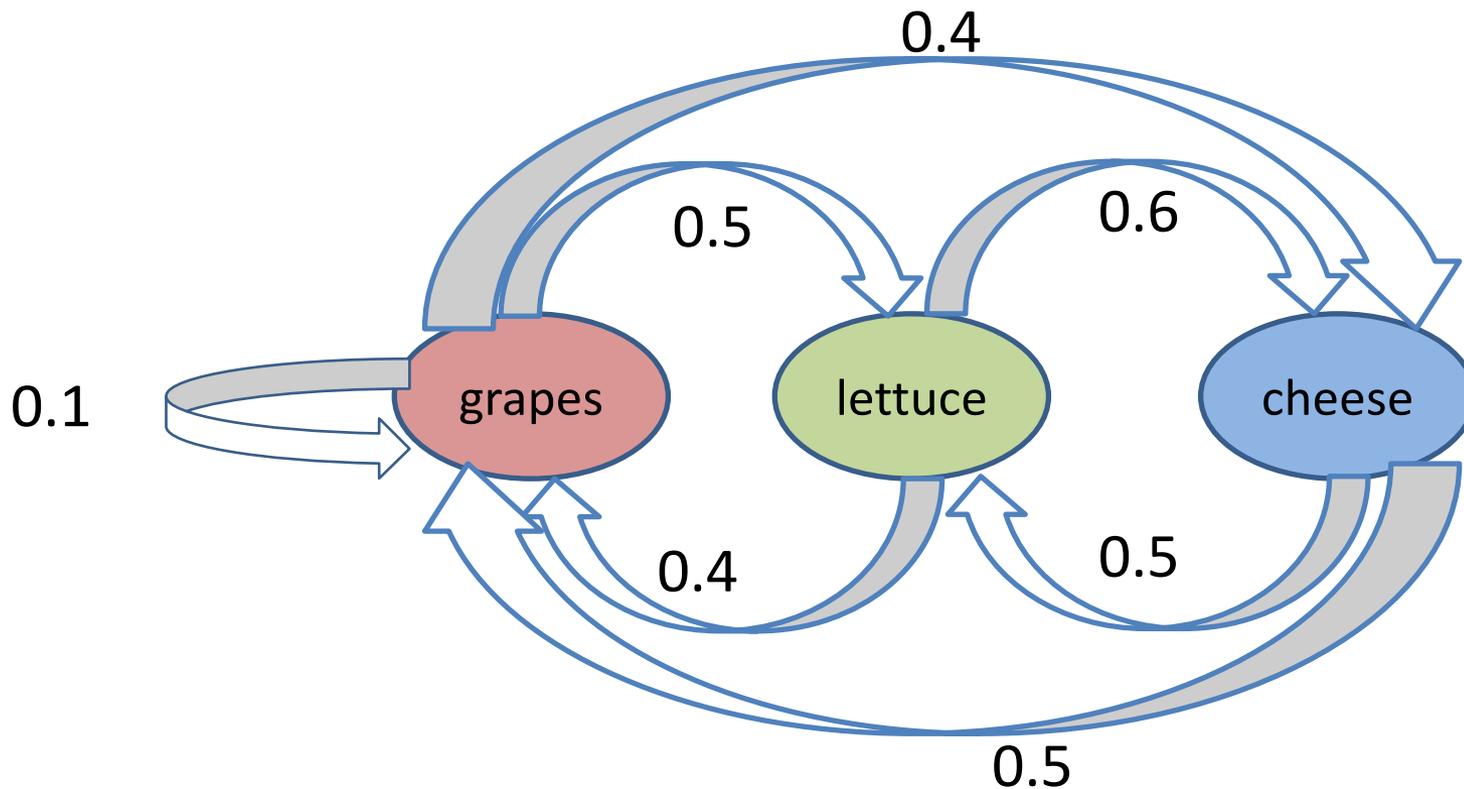
# Application of Bayes' theorem in gene regulatory network inference

$$P(M | D) = P(X) = \prod_i P(X^i | \text{Pa}(X^i))$$

- The numbers of the possible network structures (M)  $\uparrow$  with the # of nodes  $\uparrow$ .
- Search of all possible structures to find the best supported one by the data is not feasible.
- Solution: We can use **Markov chain Monte Carlo (MCMC)** simulation to identify a particular amount (e.g. 1,000) of plausible networks.
- A consensus network is obtained by combining these plausible networks.

# Markov chain

- State of the system at time “t+1” is predicted by solely using the state at time “t”, i.e. previous states (t-1, t-2,..) will not be used for predicting t+1.



# Markov chain Monte Carlo (MCMC)

- MCMC methods are a class of sampling algorithms.
- It is used for sampling from a probability distribution based on a **Markov chain** model that has the desired distribution as its equilibrium distrib.
- The state of the chain after a number of steps is used as a sample of the desired distribution.
- The quality of the sample  $\uparrow$  as the number of steps  $\uparrow$ .
- MCMC methods mainly used for numerical approximations of the integrals

# MCMC in the BN problem

- MCMC algorithm to identify 1000 of networks:
  - Start with a null network (for each 1000 cases) including prior edges.
  - Make random changes on each network such as:
    - Flip,
    - Add, and
    - Delete individual edges;
  - Accept the random changes made above that lead to an improvement in the fit of the network to the data.  $P(M | D) \approx P(D | M) \times P(M)$
  - The fitness is assessed by **Bayesian information criterion (BIC)**, which also penalizes the network if the #of the parameters (complexity)  $\uparrow$ .
  - **Note:**  $BIC = -2 \ln P(D | M) + k \ln(n)$  where  $n$  is # of data points in  $D$  (sample #);  $k = \#$  of the parameters to be estimated (parent # of node  $i$ ).
- Create a consensus network by combining above 1000 networks and check for the DAG feature of the final network.

# Comparison of BN tools

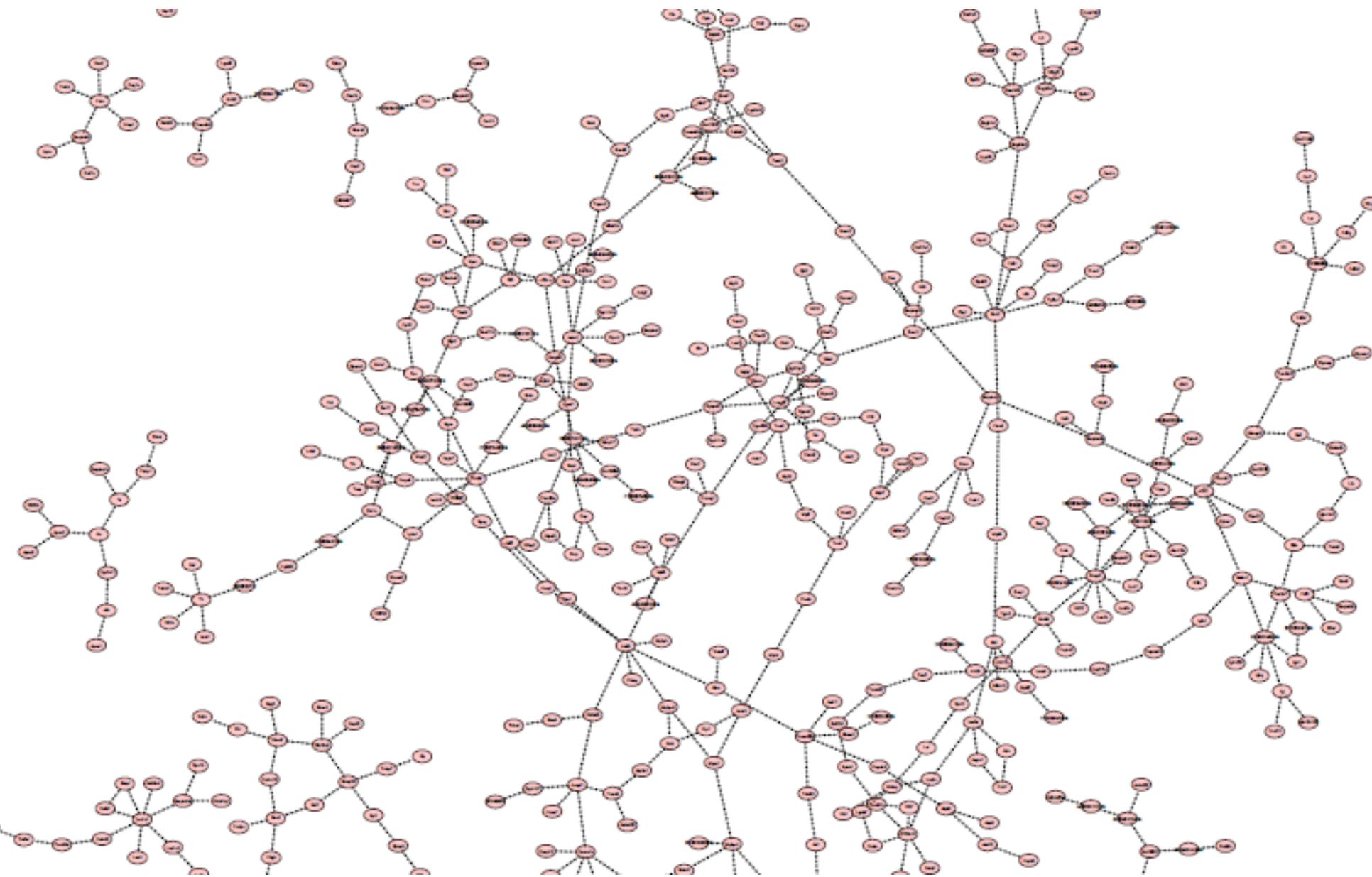
# Software packages

Tool	Platform	Scoring method(s)	Data type	Multi processing?	Prior info support?
<b>RIMBANet</b>	Perl+C	BIC	Continuous + Discrete	No	Yes
<b>Bnfinder</b>	Python	BIC , BDE, MIT	Continuous + Discrete	Yes	Yes
<b>sparsebn</b>	R	BIC+CV	Continuous + Discrete	No	<b>No</b>
<b>bnlearn</b>	R	BIC, AIC, BDE	Continuous + Discrete	Yes	Yes
<b>catnet</b>	R	BIC, AIC	Continuous + Discrete	Yes	Yes

- BIC: Bayesian information criterion
- BDE: Bayesian Dirichlet equivalence
- MIT:  $\chi^2$  distance based Mutual information test
- AIC: Akaike information criterion
- CV: cross validation

# Example - 500 genes 100 samples from an aorta tissue expression dataset

Tool	Scoring method	Runtime (sec)	Edge #
<b>RIMBANet</b>	BIC	5,586 sec (~1.30 hour) on 1 node	502
<b>Bnfinder</b>	BIC	12,710 sec (~3.30 hours) on 8 nodes (Could not finish in 24h on 1 node)	479
<b>sparsebn</b>	BIC	2,409 sec on 1 node <b>Prior info integration is not available*</b>	5,019
<b>bnlearn</b>	BIC	Could not finish in 24h (on 8 nodes)	--
<b>catnet</b>	BIC	Could not finish in 24h (on 8 nodes)	--



# Summary

- Gene networks provide us to:
  - Identify biological mechanisms and molecular subnetworks underlying common human diseases.
  - Integrating diverse type of multi-dimensional biological datasets.
  - Predicting the key driver genes in disease-related subnets.
- After presenting our data-driven findings, experimentalists can conduct *in vivo* and/or *in vitro* experiments to test our candidate genes on a given tissue, for certain hallmarks of a disease/disorder.

# Multiple sequence alignment (MSA)

# Outline

- Multiple sequence alignment (MSA)
- Introduction to MSA
- Methods of MSA
  - Progressive global alignment
  - Iterative methods
  - Alignments based on locally conserved patterns

# Motivation

- Similar genes can be conserved across species that perform similar or identical functions.
- Many genes are represented in highly conserved forms across organisms.
- By performing a simultaneous alignment of multiple sequences having similar or identical functions
  - we can gain information about
    - which regions have been subject to mutations over evolutionary time
    - and which are evolutionarily conserved.
  - Such knowledge tells
    - which regions or domains of a gene are critical to its functionality.
- Sometimes genes that are similar in sequence can be mutated or rearranged to perform an altered function.
  - By looking at multiple alignments of such sequences,
    - we can tell which changes in the sequence have caused a change in the functionality.

# Multiple sequence alignment

- a collection of three or more protein (or nucleic acid) sequences that are partially or completely aligned.
- Homologous residues are aligned in columns across the length of the sequences.
  - These aligned residues are homologous in an evolutionary sense:
    - they are presumably derived from a common ancestor.
  - The residues in each column are also presumed to be homologous in a structural sense:
    - aligned residues tend to occupy corresponding positions in the three-dimensional structure of each aligned protein.

# Multiple sequence alignment

- Multiple sequence alignment yields
  - information concerning the structure and function of proteins,
  - can help lead to the discovery of important sequence domains or motifs with biological significance
    - while at the same time uncovering evolutionary relationships among genes.
- In multiple sequence alignment, the idea is
  - to take three or more sequences, and align them
  - so that the greatest number of similar characters are aligned in the same column of the alignment.

# Multiple sequence alignment

- The difficulty with multiple sequence alignment is that
  - now there are a number of different combinations of
    - matches,
    - insertions,
    - and deletions
  - that must be considered when looking at several different sequences.
- Methods to guarantee the highest scoring alignment are not feasible.
- Therefore, approximation methods are put to use in multiple sequence alignment.

When and why are multiple sequence alignments used?

- If a protein (or gene) you are studying is related to a larger group of proteins, this group membership can often provide insight into the likely function, structure, and evolution of that protein.
- Most protein families have distantly related members.
  - Multiple sequence alignment is a far more sensitive method than pairwise alignment to detect homologs
- When the output of any database search (such as a BLAST search) is examined, a multiple sequence alignment format can be extremely useful to reveal conserved residues or motifs in the output.

When and why are multiple sequence alignments used?

- Each human genome harbors ~11,000 nonsynonymous single-nucleotide variants (causing an amino acid substitution) of which ~300 are predicted to be deleterious.
  - Algorithms that predict whether variants are harmful often rely on DNA and/or protein multiple sequence alignments to assess cross-species conservation.
    - Deleterious variants tend to occur at more conserved positions.
- Analysis of population data can provide insight into many biological questions involving evolution, structure, and function.

## When and why are multiple sequence alignments used?

- When the complete genome of any organism is sequenced, a major portion of the analysis consists of defining the protein families to which all the gene products belong.
- Database searches effectively perform multiple sequence alignments, allowing comparisons of each novel protein (or gene) to the families of all other known genes.
- The most critical part of making a phylogenetic tree is to produce an optimal multiple sequence alignment.
- The regulatory regions of many genes contain consensus sequences for transcription factor-binding sites and other conserved elements.
  - Many such regions are identified based on conserved noncoding sequences that are detected using multiple sequence alignment.

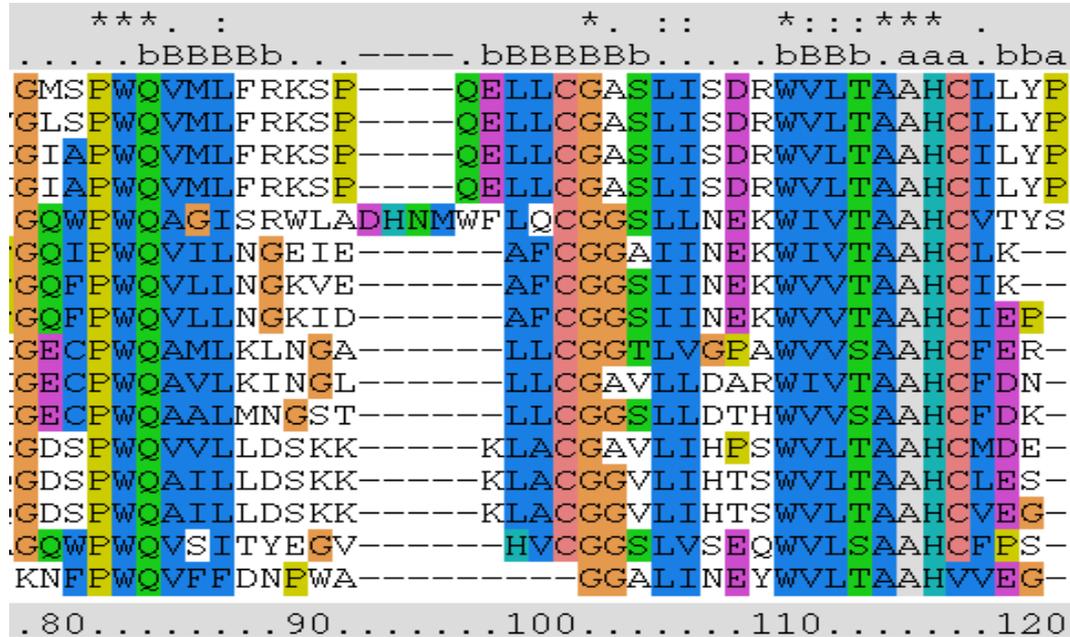
# Example Multiple Alignment

	1	1	10	20	30
<b>NONAME</b>	1	VSLTCL	VKGFYPSD	IAVEWESNG	--
<b>NONAME#2</b>	1	VTISCTGT	SSNIGS	ITVNWYQQLPG	
<b>NONAME#8</b>	1	VTISCTGSS	SNIGAG	NHVKWYQQLPG	
<b>NONAME#3</b>	1	LRLSCS	SSGFIFSS	YAMYWVRQAPG	
<b>NONAME#4</b>	1	LSLTCT	VSGTSFDD	YYSTWVRQPPG	
<b>NONAME#5</b>	1	PEVTCVVVD	VSHEDPQVKFNWYVDG	--	
<b>NONAME#6</b>	1	ATLVCL	ISDFYPGA	VTVAWKADS	--
<b>NONAME#7</b>	1	AALGCL	VKDYFPEP	VTVSWNSG	---
<b>Consensus</b>	1	VTL SCT	VS F S	V V W	Q PG

Ready positives: 59.3% identity: 7.4% al

- Example multiple alignment of 8 immunoglobulin sequences.

# Example MSA

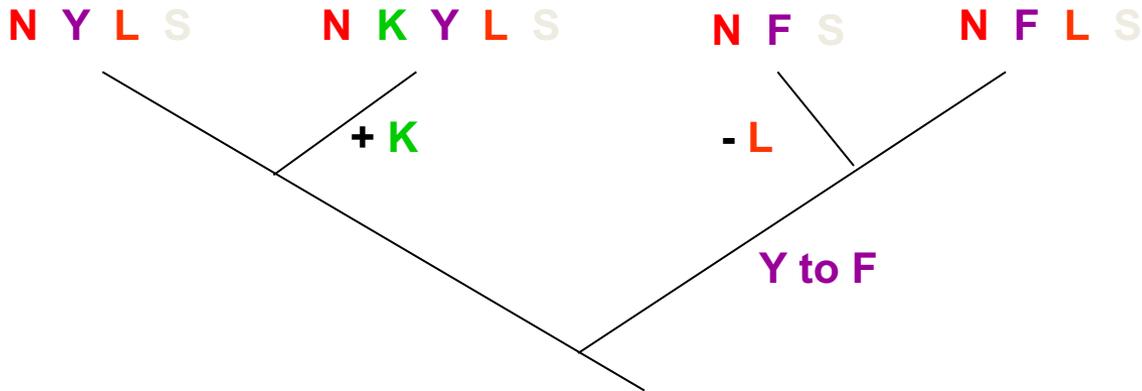


- Each row is a different protein sequence
- Each column is a different aligned position

# Relationship of MSA to Phylogenetic analysis

once the msa has been found, the number or types of changes in the aligned sequences may be used for a phylogenetic analysis

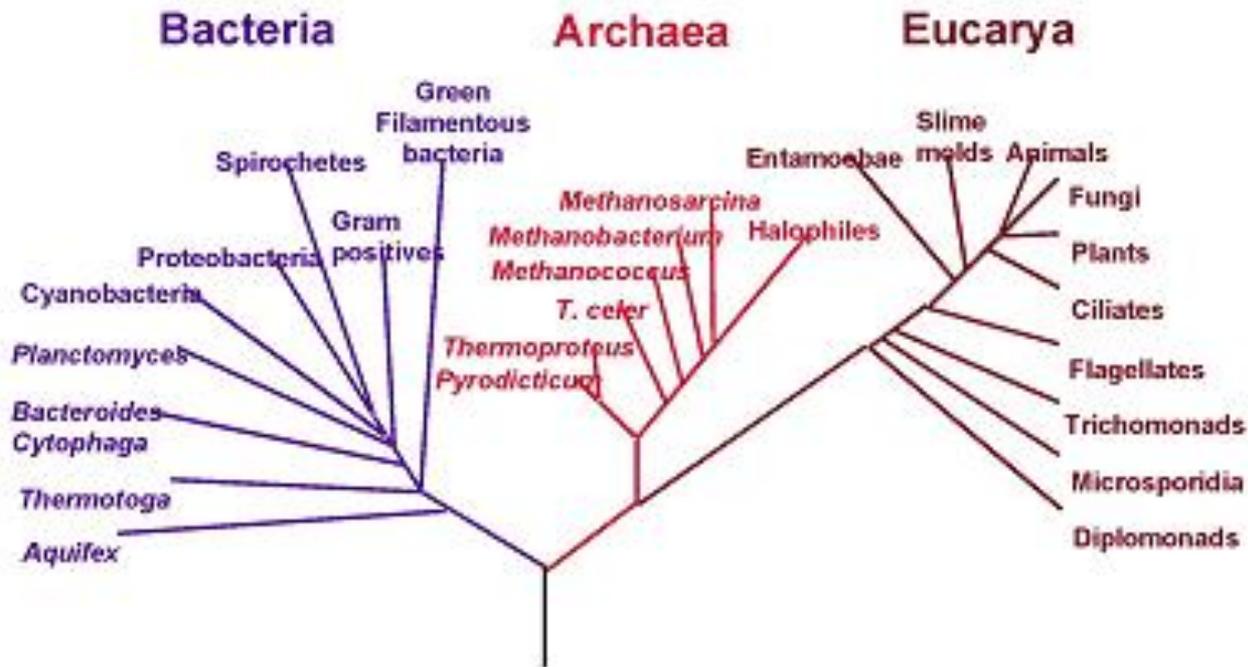
seqA	N	-	F	L	S
seqB	N	-	F	-	S
seqC	N	K	Y	L	S
seqD	N	-	Y	L	S



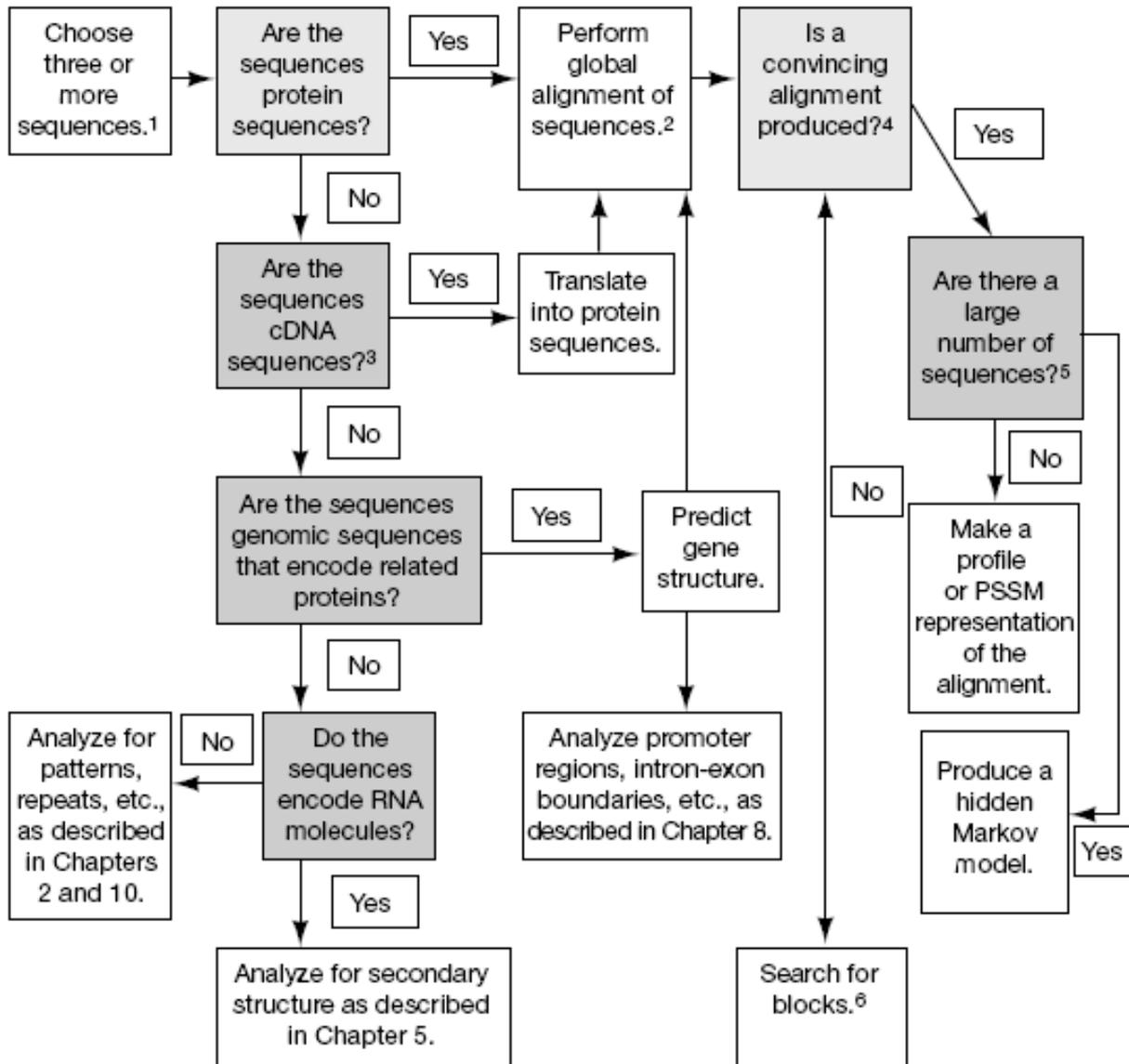
hypothetical evolutionary tree that could have generated three sequence changes

# Phylogenetic analysis

## Phylogenetic Tree of Life



# M A C A



## Approaches to Multiple Alignment

- There are many approaches to multiple sequence alignment;
  - in the past decade many dozens of programs have been introduced
    - <https://doi.org/10.1093/bib%2F6.1.6>
- We will consider four approaches to multiple sequence alignment:
  - Exact methods (Dynamic Programming)
  - Progressive Alignment
  - Iterative Alignment
  - Statistical Modeling

## Dynamic Programming Approach

- Exact methods of multiple alignment use **dynamic programming**
  - guaranteed to find optimal solutions.
- Dynamic programming with two sequences
  - Relatively easy to code
  - Guaranteed to obtain optimal alignment
- Can this be extended to multiple sequences?

## Dynamic Programming With 3 Sequences

- Consider the following amino acid sequences to align

VSNS, SNA, AS

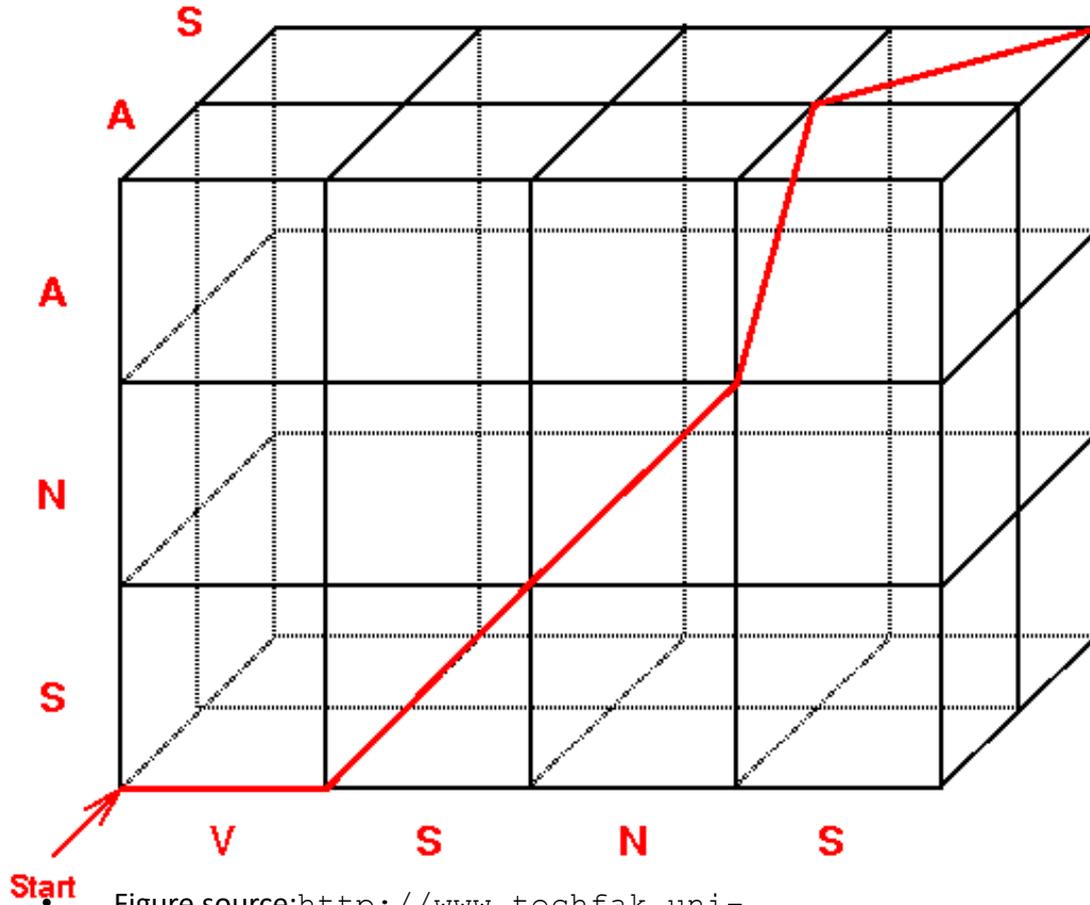
- Instead of filling a two dimensional matrix as we did with two sequences,
  - we now fill a three dimensional space.
- Put one sequence per axis (x, y, z)
- Three dimensional structure results

# Dynamic Programming With 3 Sequences

## Possibilities:

- All three match;
- A & B match with gap in C
- A & C match with gap in B
- B & C match with gap in A
- A with gap in B & C
- B with gap in A & C
- C with gap in A & B

# Dynamic Programming With 3 Sequences



V S N \_ S  
 \_ S N A \_  
 \_ \_ \_ A S

Figure source: <http://www.techfak.uni-bielefeld.de/bcd/Curric/MulAli/node2.html#SECTION00020000000000000000>

# Dynamic Programming

- Suppose the length of each sequence is  $n$  residues.
- If there are two such sequences,
  - then the number of comparisons needed to fill in the scoring matrix
    - is  $n^2$ ,
    - since it is a two-dimensional matrix.
- The number of comparisons needed to fill in the scoring cube
  - when three sequences are aligned
    - is  $n^3$ ,
  - and when four sequences are aligned,
    - is  $n^4$ .

- Thus, as the number of sequences increases,
  - the number of comparisons needed increases exponentially, i.e.  $n^N$ 
    - where  $n$  is the length of the sequences,
    - and  $N$  is the number of sequences.
- Thus, without any changes to the dynamic programming approach, this becomes impractical for even a small number of short sequences rather quickly.

# Example

**2 protein sequences length = 300, excluding gaps**

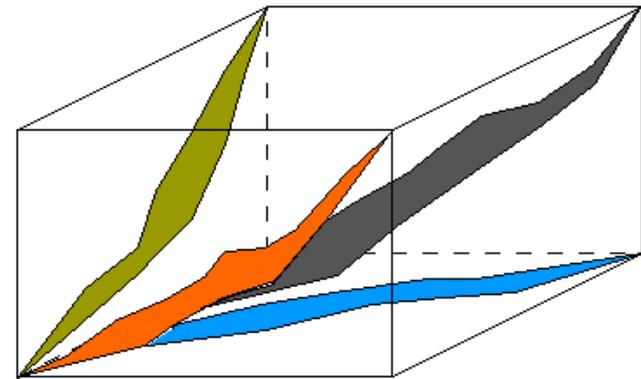
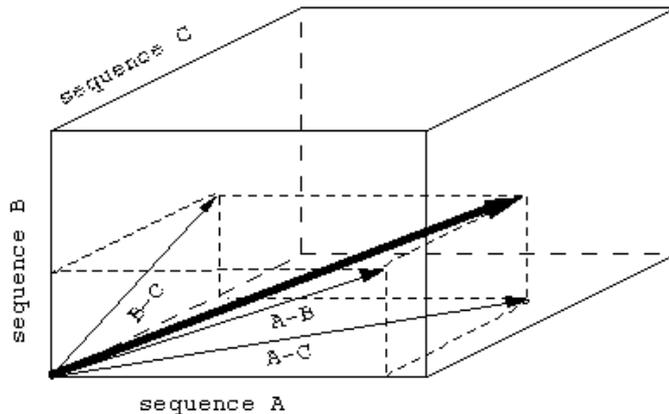
number of comparisons by  
dynamic programming  $300^2 = 9 * 10^4$

**3 protein sequences length = 300, excluding gaps**

number of comparisons  
by dynamic programming  $300^3 = 2.7 * 10^7$

# Reduction of space and time

- Carrillo and Lipman:
  - multiple sequence alignment space bounded by pairwise alignments
  - Projections of these alignments lead to a bounded alignments
    - Following figures show how the two dimensional search spaces can be projected into a three dimensional volume that can be searched



# Reduction of space and time

- Idea for reduction of memory and computations:
  - Multiple sequence alignment imposes an alignment on each of the pairs of sequences.
- Alignments found for each of the pairs of sequences can impose bounds on the location of the MSA within the cube (three sequences) or N-dimensional space (N sequences).
  - Step 1
    - Find pairwise alignment for sequences.
  - Step 2
    - Trial *msa* produced by predicting a phylogenetic tree for the sequences
  - Step 3
    - Sequences multiply aligned in the order of their relationship on the tree

# Reduction of space and time

- This is a heuristic alignment
- Therefore the alignment is not guaranteed to be optimal
- Alignment provides a limit to the volume within which optimal alignments are likely to be found

# MSA

- MSA: Developed by Lipman, 1989
- Incorporates extended dynamic programming
- MSA calculates the multiple alignment score within the lattice by adding the scores of the corresponding pairwise alignments in the multiple sequence alignment.
- This measure is known as the **sum of pairs (SP)** measure.
- The optimal alignment is based on the best SP score.

# Scoring of msa's

- MSA uses **Sum of Pairs (SP)**
  - Scores of pair-wise alignments in each column added together
  - Columns can be weighted to reduce influence of closely related sequences
  - Weight is determined by distance in phylogenetic tree

# Sum of Pairs Method

- The **sum of pairs** method scores all possible combinations of pairs of residues in a column of a multiple sequence alignment.
- For instance, consider the alignment

E    C    S    Q                    (1)

S    N    S    G                    (2)

S    W    K    N                    (3)

S    C    S    N                    (4)

- Since there are four sequences,
  - there will be six different alignments to consider for each column.
- The alignments, listed by the sequence number are listed as follows:

(1)-(2); (1)-(3); (1)-(4); (2)-(3); (2)-(4); (3)-(4)



- Problem with this approach:
  - more closely related sequences will have a higher weight
- The MSA program gets around this by calculating weights to associate to each sequence alignment pair.
- The weights are assigned based on the predicted tree of the aligned sequences.

# Summary of MSA

1. Calculate all pairwise alignment scores
2. Use the scores to predict tree
3. Calculate pair weights based on the tree
4. Produce a heuristic msa based on the tree
5. Calculate the maximum weight for each sequence pair
6. Determine the spatial positions that must be calculated to obtain the optimal alignment
7. Perform the optimal alignment
8. Report the weight found compared to the maximum weight previously found

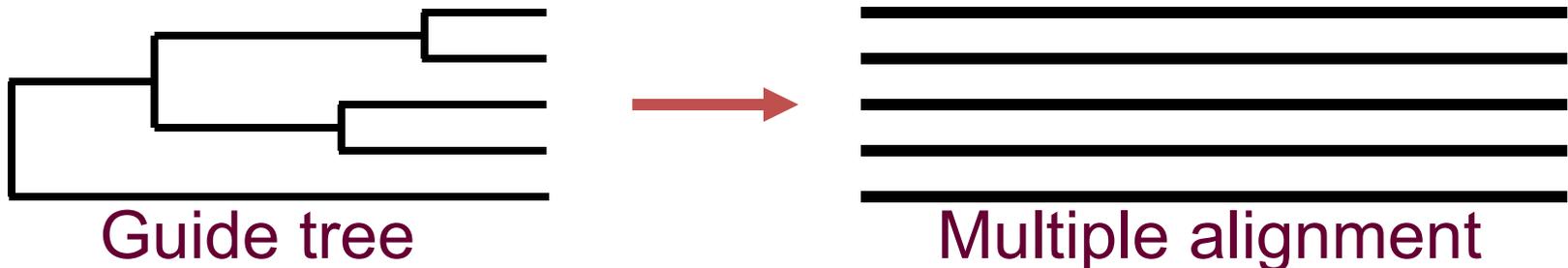
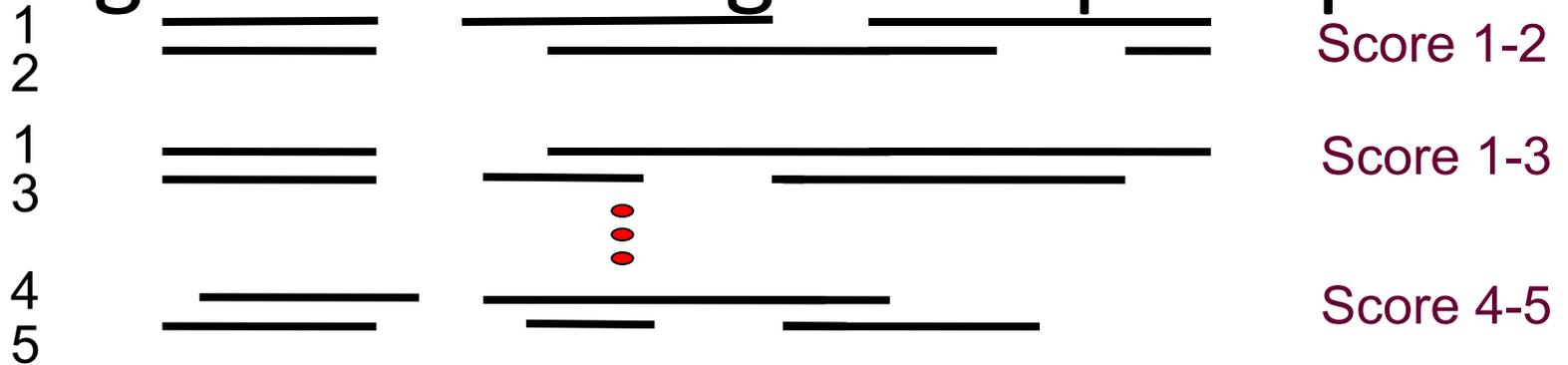
# Progressive Alignments

- MSA program is limited in size
- The approach of progressive alignment is
  - to begin with an alignment of the most alike sequences,
  - and then build upon the alignment using other sequences.
- Progressive alignments work by
  - first aligning the most alike sequences using dynamic programming,
  - and then progressively adding less related sequences to the initial alignment.

## Progressive multiple sequence alignment

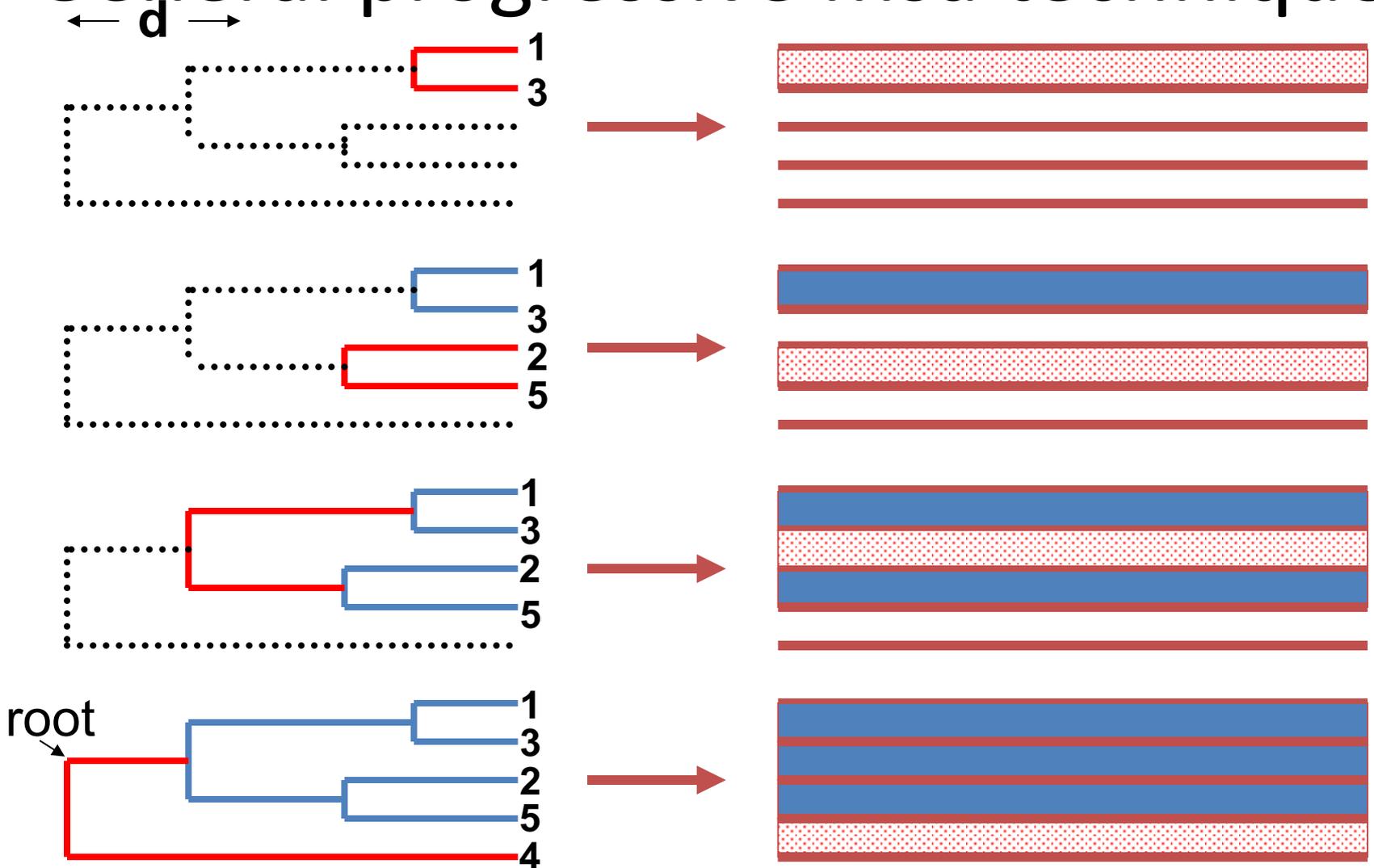
- alignment on each of the pairs of sequences
- next, trail msa is produced by first predicting a phylogenetic tree for the sequences
- sequences are then multiply aligned in order of their relationship on the tree
  - starting with the most related sequences
  - then progressively adding less related sequences to the initial alignment
- used by PILEUP and CLUSTALW
- not guaranteed to be optimal

# Progressive msa - general principles



# General progressive msa technique

(follow generated tree)



# CLUSTALW and CLUSTALX

- CLUSTALW and CLUSTALX are progressive alignment programs that follow the following steps:
  - Perform pairwise alignments of all of the sequences
  - Use the alignment scores to produce a phylogenetic tree using neighbor-joining methods
  - Align the sequences sequentially, guided by the phylogenetic relationships indicated by the tree

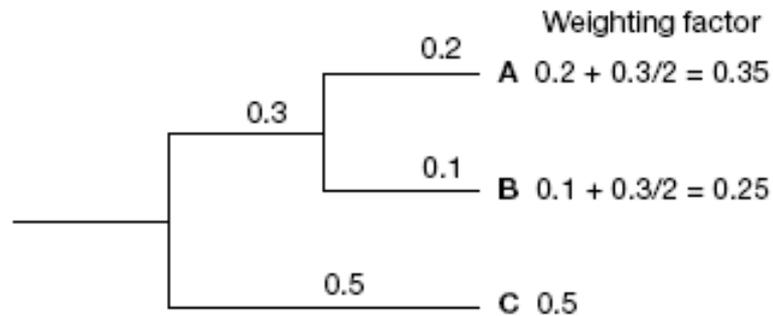
# CLUSTALW

- The initial pairwise alignments are calculated using an enhanced dynamic programming algorithm, and
- the genetic distances used to create the phylogenetic tree are calculated by dividing the total number of mismatched positions by the total number of matched positions.

# CLUSTALW

- Alignments are associated a weight based on their distance from the root node (next slide)
- Gaps are added to an existing profile in progressive methods
- CLUSTALW incorporates a statistical model in order to place gaps where they are most likely to occur

**A. Calculation of sequence weights**



**B. Use of sequence weights**

Column in alignment 1

Sequence A (weight a)    .....K.....

Sequence B (weight b)    .....I.....

Column in alignment 2

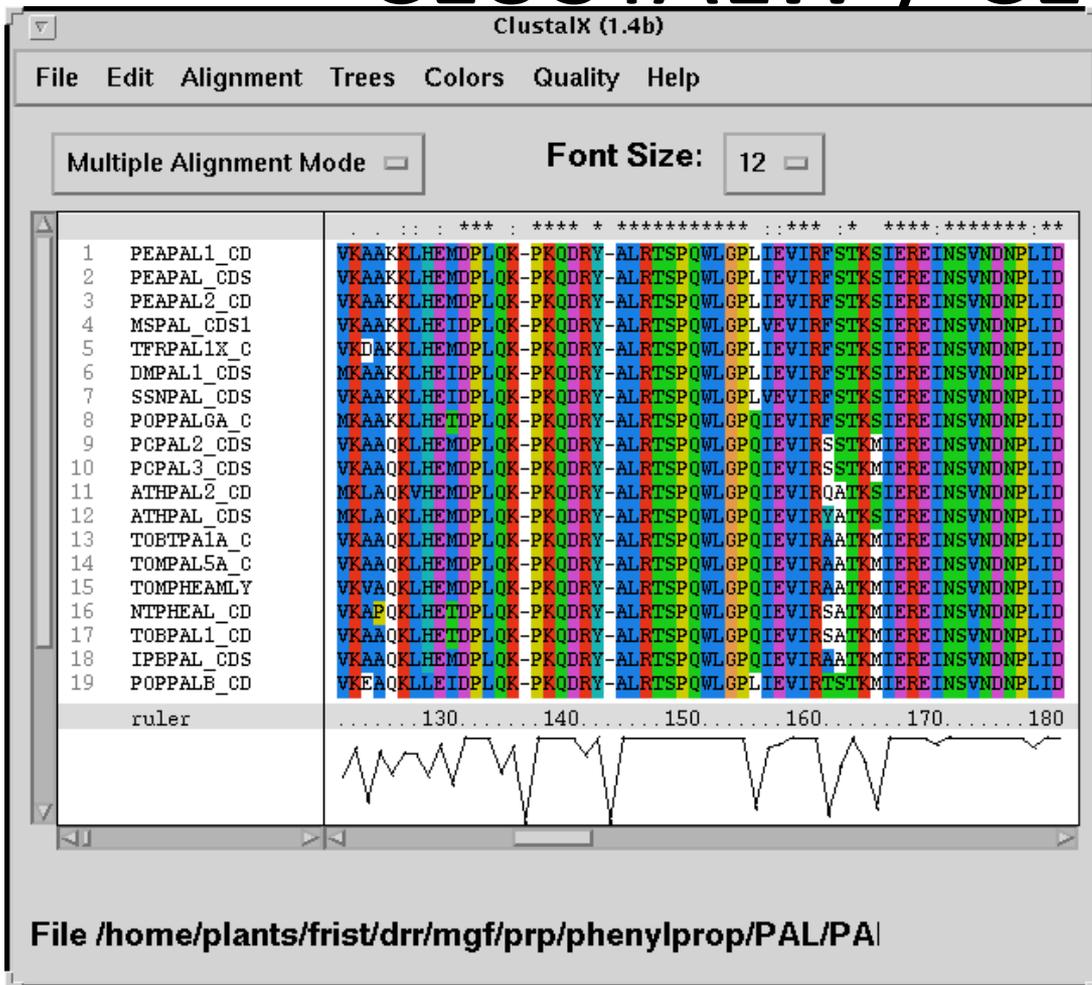
Sequence C (weight c)    .....L.....

Sequence D (weight d)    .....V.....

Score for matching these two column in an msa =

$$[ a \times c \times \text{score}(K,L) + a \times d \times \text{score}(K,V) + b \times c \times \text{score}(I,L) + b \times d \times \text{score}(I,V) ] / 4$$

# CLUSTALW / CLUSTALX



- ‘W’ stands for “weighting”
  - ability to provide weights to sequence and program parameters
- CLUSTALX – with graphical interface
- provides global msa
- Not constructed to perform local alignments.
- Similarity in small regions is a problem.
- Problems with large insertions.
- Problems with repetitive elements, such as domains.
- ClustalW does not guarantee an optimal solution

# CLUSTALW

- <http://www.ebi.ac.uk/clustalw/>

# PILEUP

- PILEUP is the multiple sequence alignment program that is part of the Genetics Computer Group (GCG) package developed at the University of Wisconsin.
- Sequences initially aligned in a pair-wise fashion using Needleman-Wunsch algorithm.
- Scores used to produce tree using unweighted pair group method using arithmetic averages (UPGMA)
- The resulting tree is then used to guide the alignment of the most closely related sequences and groups of sequences

# PILEUP

- very similar to CLUSTALW
- part of the genetic computer group (GCG)
- does not guarantee optimal alignment
- plots a cluster dendrogram of similarities between sequences



## Shortcoming of Progressive Approach

- Dependence upon initial pair-wise sequence alignments
  - Ok if sequences are similar
  - Errors in alignment propagated if not similar
- Suitable scoring matrices and gap penalties must be chosen to apply to the sequences as a set

# Iterative Methods

- Iterative alignment methods begin by making an initial alignment of the sequences.
- These alignments are then revised to give a more reasonable result.
- The objective of this approach is to improve the overall alignment score
- Alignment is repeatedly refined
- Selection of groups is based on the phylogenetic tree
- Programs using iterative methods:
  - MultAling
  - PRRP
  - DIALIGN

# MultAlign

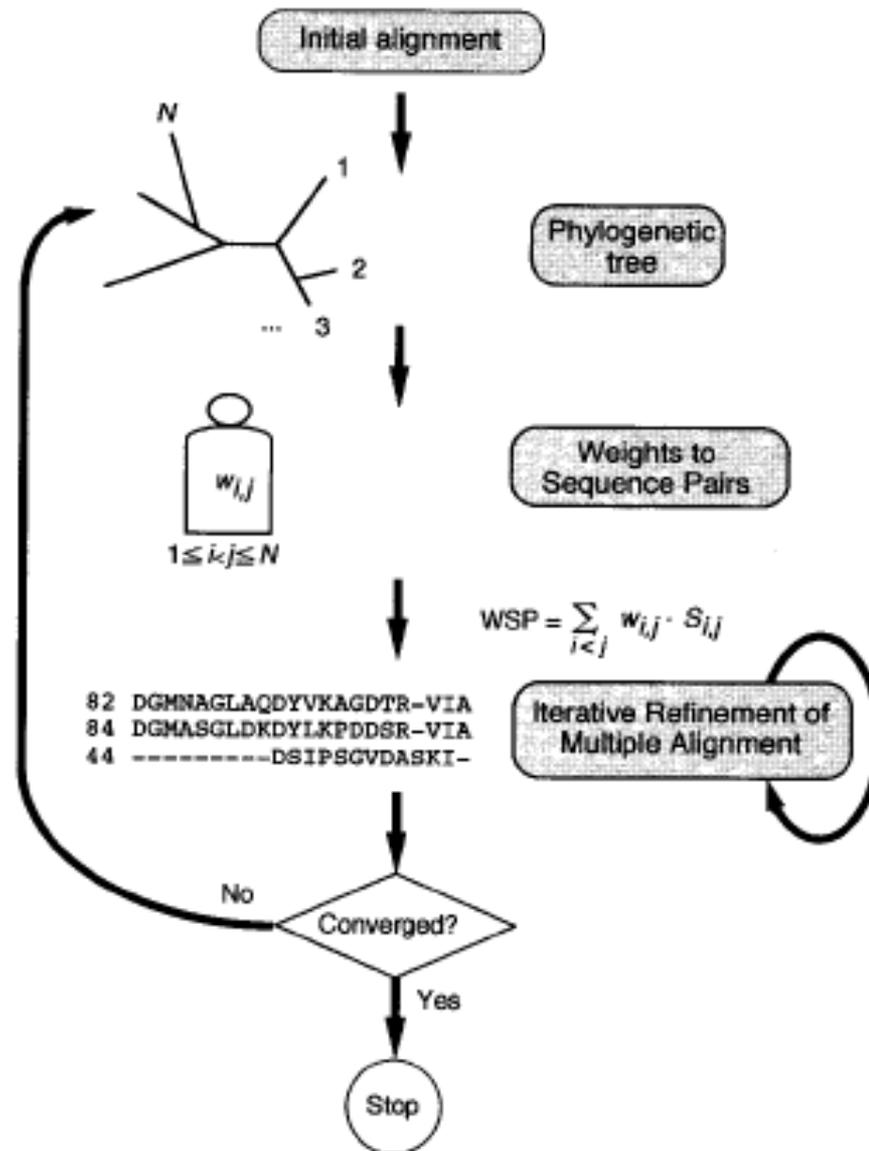
- Pairwise scores recalculated during progressive alignment
- Tree is recalculated
- Alignment is refined

# PRRP

- Initial pairwise alignment predicts tree
- Tree produces weights
- Locally aligned regions considered to produce new alignment and tree
- Continue until alignments converge

Iter

SA



# DIALIGN

- Pairs of sequences aligned to locate ungapped aligned regions
- Diagonals of various lengths identified
- Collection of weighted diagonals provide alignment

# Genetic Algorithms

- The goal of genetic algorithms used in sequence alignment is to generate as many different multiple sequence alignments by rearrangements that simulate gaps and genetic recombination events.
- SAGA (Serial Alignment by Genetic Algorithm) is one such approach that yields very promising results, but becomes slow when more than 20 sequences are used.

# Genetic Algorithm Approach

- 1) Sequences (up to 20) written in row, allowing for overlaps of random length – ends padded with gaps (100 or so alignments)

```
XXXXXXXXXX-----  
-----XXXXXXXXXX  
--XXXXXXXXXX-----
```

# Genetic Algorithm Approach

- 2) The initial alignments are scored by the sum of pairs method.
  - Standard amino acid scoring matrices and gap open, gap extension penalties are used
- 3) Initial alignments are replaced to give another generation of multiple sequence alignments
  - One half of the multiple sequence alignments are chosen to proceed to the next generation unchanged (natural selection).
    - This half is chosen by assigning probabilities to each sequence based on an inverse proportion of their SP scores (the best alignments, since the SP scores are weighted according to their distance from the parent).
  - The other half of the alignments are sent to the next generation, but are first subject to mutation.

# Genetic Algorithm Approach

## 4) Mutation

- In the mutation process, gaps are inserted into the sequences subject to mutation and rearranged in an attempt to create a better scoring alignment.
- In this step
  - the sequences subject to mutation split into two sets based on estimated phylogenetic tree
  - gaps of random lengths inserted into random positions in the alignment

# Genetic Algorithm Approach

- Mutations:

- XXXXXXXXX      XXX---XXX—XX
- XXXXXXXXX      XXX---XXX—XX
- XXXXXXXXX      X—XXX---XXXX
- XXXXXXXXX      X—XXX---XXXX
- XXXXXXXXX      X—XXX---XXXX

# Genetic Algorithm Approach

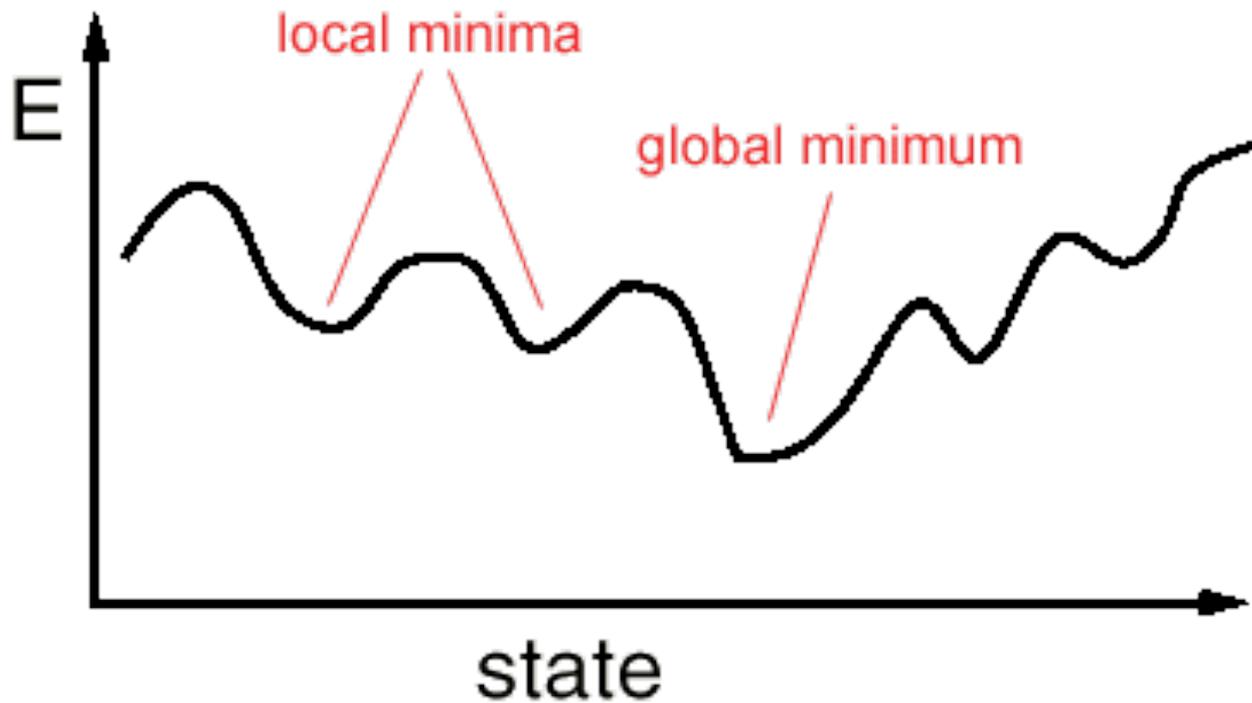
- 5) Recombination of two parents to produce next generation alignment is accomplished
- 6) The next generation is evaluated going back to step 2, and steps 2-5 are repeated a number (100-1000) times.
  - The best scoring multiple sequence alignment is then obtained (note that it may not be the optimal scoring alignment).
- 7) The entire process is repeated several times, starting from a different initial alignment each time.
  - The best scoring multiple sequence alignment is then chosen and reported to the user.

# Simulated Annealing

- Another approach to sequence alignment that works in a manner similar to genetic algorithms is simulated annealing.
- In these approaches, you begin with a heuristically determined multiple sequence alignment that is then changed using probabilistic models that identifies changes in the alignment that increase the alignment score.
- The drawback of simulated annealing approaches is that you can get stuck finding only the locally optimal alignment rather than the alignment score that is globally optimal
- Rearranges current alignment using probabilistic approach to identify changes that increase alignment score

# Simulated Annealing

## Simulated Annealing



<http://www.cs.berkeley.edu/~amd/CS294S97/notes/day15/day15.html>

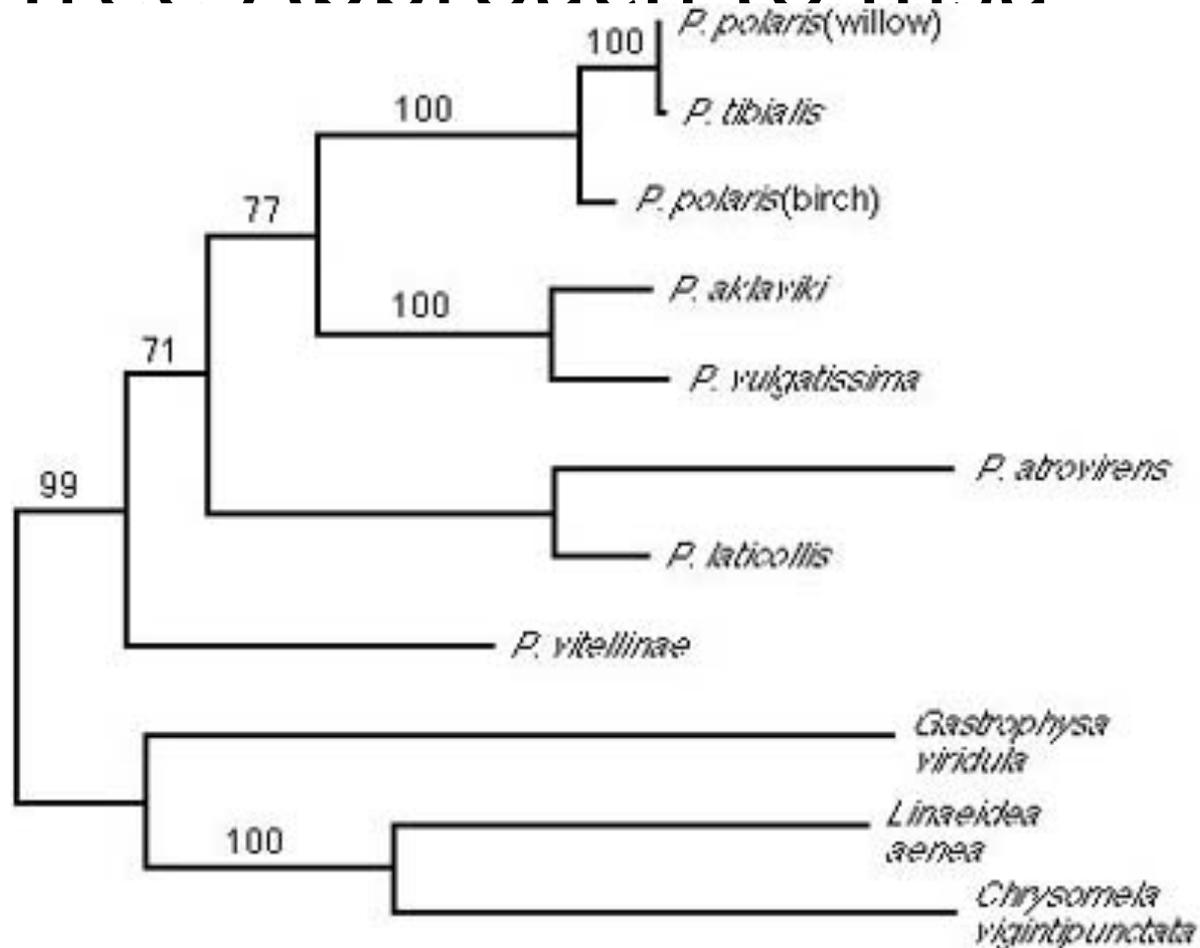
# Group Approach

- Sequences aligned into similar groups
- Consensus of group is created
- Alignments between groups is formed
  
- **EXAMPLES: PIMA, MULTAL**

# Tree Approach

- Tree created
- Two closest sequences aligned
- Consensus aligned with next best sequence or group of sequences
- Proceed until all sequences are aligned

# Tree Approach to msa



- [www.sonoma.edu/users/r/rank/research/evolhost3.html](http://www.sonoma.edu/users/r/rank/research/evolhost3.html)

# Tree Approach to msa

- PILEUP, CLUSTALW and ALIGN
- TREEALIGN rearranges the tree as sequences are added, to produce a maximum parsimony tree (fewest evolutionary changes)

# Profile Analysis

- Create multiple sequence alignment
- Select conserved regions
- Create a matrix to store information about alignment
  - One row for each position in alignment
  - one column for each residue; gap open; gap extend

# Profile Analysis

- Profile can be used to search target sequence on database for occurrence
- Drawback: profile is skewed towards training data

# **Local Multiple Sequence Alignment Sequence File Formats**

# Localized Alignments

- Just like with pairwise alignments, we may not be interested in the global alignment of multiple sequences, but rather only specific regions that are conserved.
- Local Alignment of **msas** are important:
  - Given regions of genomic DNA occurring upstream or before a certain gene, there might be sequences where transcription factors bind to the DNA so that the gene can be transcribed.
  - Thus, if we are interested in determining if there is any signal in the regions upstream of a certain family of genes across several different organisms, it would be important to only find the conserved region, and not try to align all of the genomic DNA
  - Localized alignments of protein sequences can yield information about conserved domains found in otherwise unrelated proteins.

# Approaches to Local Alignment

- Profile Analysis
- Block Analysis
- Pattern-searching or statistical methods

# Profile Analysis

- Profiles are found by first multiply aligning the sequences, determining which regions are the most highly conserved, and
- then creating a scoring matrix for the alignment of the highly conserved region.
- Profile is composed of:
  - Columns: one for each residue;
    - columns for insertions and deletions as well
  - Rows: one for each position in the conserved region or motif

# Profile Analysis

- Profiles describe a msa by a scoring matrix:

	1 V	2 L	3 V	4 K	5 A	6 G	7 S	8 S	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W
1WGVL	V	3	-2	3	4	0	4	-1	3	-1	4	4	1	1	1	-2	1	2	6	-6	.	.	.	.	.	.	.
2LLSP	L	2	-2	-1	3	0	-1	3	-1	6	5	-1	3	0	-1	3	1	4	1	.	.	.	.	.	.	.	.
3VVVV	V	2	2	-2	-2	2	2	-3	1	-2	8	6	-2	1	-2	-2	0	2	15	-9	.	.	.	.	.	.	.
4KEAT	K	6	-2	5	6	-5	4	1	0	5	-2	0	3	3	3	1	3	6	0	-6	.	.	.	.	.	.	.
5APLP	P	6	-1	0	1	-2	2	0	1	0	2	2	0	8	2	0	2	2	3	-5	.	.	.	.	.	.	.
6GGGG	G	7	1	7	5	-6	15	-1	-3	0	-4	-3	4	3	6	1	6	2	-1	-6	.	.	.	.	.	.	.
7SSQE	S	4	-1	7	7	-6	7	2	-2	2	-3	-2	4	3	6	1	6	2	-1	-6	.	.	.	.	.	.	.
8SSSTP	S	4	4	2	2	-4	4	-1	0	2	-3	-2	2	7	0	1	10	6	0	-2	.	.	.	.	.	.	.

# Profile Searches

Once a profile is created, it can be used to search a target sequence or database for possible matches to the profile using the profiles scores to evaluate the likelihood at each position.

Profile scores evaluate likelihood of a match at each position

# Drawback to Profiles

- Profiles only as representative as the variation in the training sets.
- Thus, there is a bias in the profile towards the training data.
- Training sets can be erroneous if not carefully constructed

# Calculating Profiles

- Each cell is the log-odds score
  - The value of an individual cell is calculated as the log odds score of finding a particular residue in a particular location in an alignment divided by the probability of aligning the two amino acids by random chance using a particular scoring scheme (such as PAM250, BLOSUM80, ...).
    - PAM (Percent Accepted Mutation)
    - BLOSUM (Blocks Substitution Matrix)
  - Additional penalties must be calculated for gap opening and gap extension in the profile as well.
- Some methods take in sequence weights as well

# Shannon Entropy

- One method to calculate the observed column variation given the expected variation in the evolutionary model is to use an information measure known as **entropy**.
  - **Entropy** is the amount of information of the observed column variation if expected variation in the evolutionary model is known
- The smaller the entropy, the more conserved a column is.

# Entropy

- The entropy ( $H$ ) for a single column is calculated by the following formula:

$$H = - \sum_{\text{residues}(a)} f_a \log(p_a)$$

- $a$ : is a residue (amino acid),
- $f_a$ : frequency of residue  $a$  in a column,
- $p_a$  : probability (expected frequency) of residue  $a$  in that column

# Entropy

- $H$  is calculated for each 20 ancestor amino acids and for a large number of evolutionary distances (PAM1, PAM2, PAM4, ...).
- The distance that gives the minimum value for  $H$  for each column-possible ancestor combination is the best estimate of the distance that generates the column diversity from that ancestor.
- This analysis provides 20 possible models ( $M_a$  for  $a = 1, 2, 3, \dots, 20$ ) as to show how the amino acid frequencies in a column ( $F$ ) may have originated.

# Entropy

- The next step in the evolutionary profile construction determines the extent to which each  $M_a$  predicts  $F$  by the Bayes conditional probability analysis.

$$P(M_a|F) = P(M_a) \times P(F|M_a) / \sum_{\text{all } a@s} P(M_a) \times P(F|M_a)$$

where the prior distribution  $P(M_a)$  is given by the background amino acid frequencies and

$$P(F|M_a) = P_{aa1}^{faa1} \times P_{aa2}^{faa2} \times P_{aa3}^{faa3} \dots \dots P_{aa20}^{faa20}$$

i.e., the product of the expected amino acid frequencies in  $M_a$  raised to the power of the fraction observed for each amino acid in the msa column.

# Entropy

- From  $P(M_a | F)$ , the weights for each of the 20 possible distributions that give rise to the msa column diversity are calculated as follows:

$$W_a = P(M_a | F) - P(M_{random} | F)$$

where  $W_a$  is the weight given to  $M_a$  and  $P(M_{random} | F)$  is calculated as above using amino acid distribution.

# Log-odds score

- Another measure of creating a profile is by using **log-odds score**.
- In this method, the  **$\log_2$  of the ratio of observed/background frequencies** is calculated for each position.
- What results is the amount of information available in an alignment given in bits.
- A new sequence can then be searched to see if it possibly contains the motif.
- Profiles can also indicate log-odds score:
  - $\text{Log}_2(\text{observed}:\text{expected})$
- Result is a bit score

# Log-odds score

- The log odds scores for the profile ( $\text{Profile}_{ij}$ ) are given by

$$\text{Profile}_{ij} = \log \left[ \sum_{\text{all } a \in \mathcal{S}} (W_{ai} \times P_{aij}) / P_{\text{random}j} \right]$$

where  $W_{ai}$  is the weight of an ancestral amino acid  $a$  at row  $i$  in the profile,  $P_{aij}$  is the frequency of amino acid distribution that best matches at row  $i$ , and  $P_{\text{random}j}$  is the background frequency of amino acid  $j$ .

# BLOCKS

- Blocks are similar to profiles in the sense that
  - they represent locally conserved regions within a multiple sequence alignment.
- However, the difference is that ...
  - blocks lack insert and delete (**indels**) positions in the sequences.
  - Instead, every column includes only matches and mismatches
- Blocks can be determined either
  - by performing a multiple sequence alignment, or
  - by searching a database for similar sequences of the same length.

# BLOCKS

- Generally determined by performing multiple alignment first
- Ungapped regions are then separated into blocks
- Algorithms have been developed for searching for blocks

# BLOCKS

- Statistical approaches to finding the most alike sequences have been proposed, such as
  - the Expectation-Maximization algorithms and
  - the Gibbs sampler.
- In any case, once a set of blocks has been determined, the information contained within the block alignment can be displayed as a sequence profile.

# BLOCKS Programs

- A global sequence alignment will usually contain ungapped regions that are aligned between multiple sequences.
- These regions can be extracted to produce blocks.
- Two widely used programs:
  - BLOCKS
  - eMOTIF

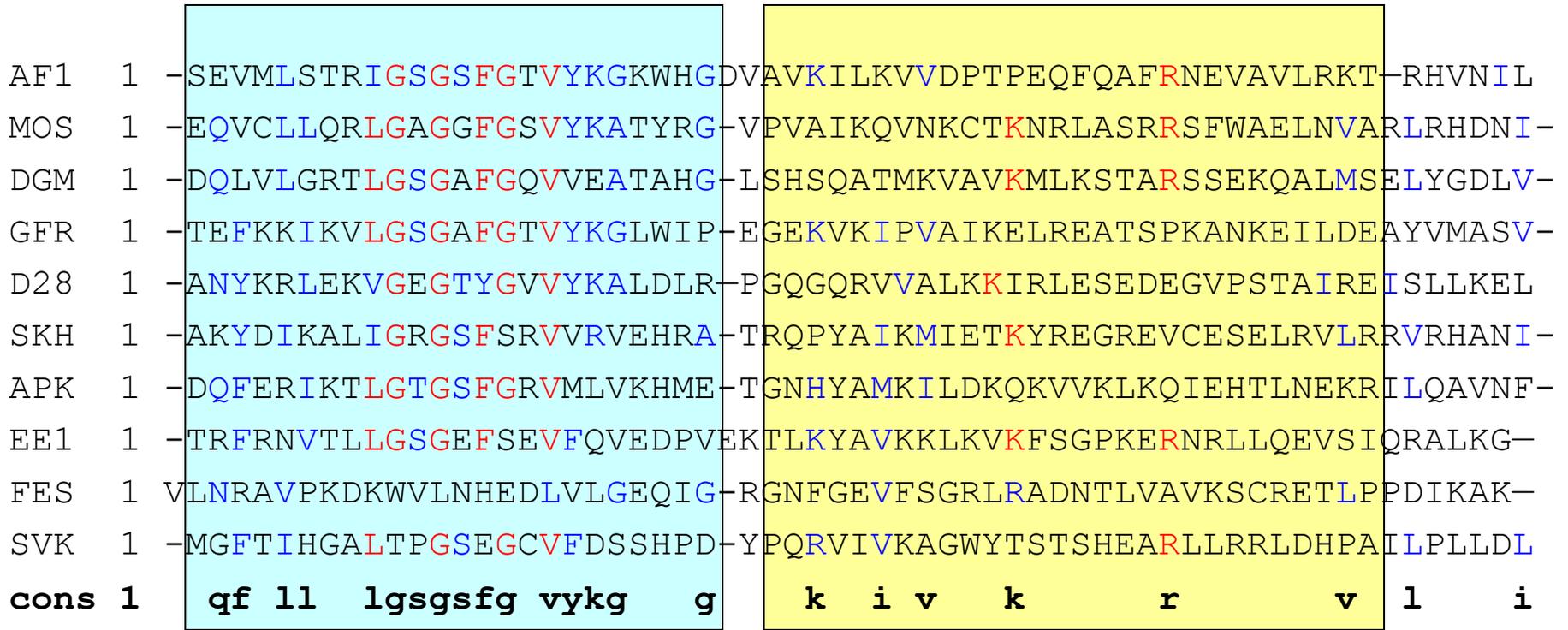
[http://www.blocks.fhcrc.org/blocks/process\\_blocks.html](http://www.blocks.fhcrc.org/blocks/process_blocks.html)

<http://dna.stanford.edu/emotif/>

- Example
  - 10 Truncated Kinase proteins
  - Approximately 75 residues in length

**>D28 CD28 S. CEREVISIAE CELL CYCLE CONTROL PROTEIN KINASE**  
 ANYKRLEKVGEGTYGVVYKALDLRPGQGQRVVALKKIRLESEDEGVPSTAIREISLLKEL  
**>SKH SKH HELA MYSTERY PUTATIVE PROTEIN KINASE**  
 AKYDIKALIGRGSFSRVVRVEHRATRQPYAIKMIETKYREGREVCESLRVLRVRHANI  
**>APK CAPK BOVINE CARDIAC MUSCLE CYCLIC AMP-DEPENDENT (ALPHA)**  
 DQFERIKTLGTGSFGRVMLVKHMETGNHYAMKILDKQKVVKLKQIEHTLNEKRILQAVNF  
**>EE1 WEE1 S. POMBE MITOTIC INHIBITOR**  
 TRFRNVTLGSGEFSEVFQVEDPVEKTLKYAVKKLKVKFSGPKERNLLQEVS IQRALKG  
**>GFR EGFR HUMAN EPIDERMAL GROWTH FACTOR RECEPTOR**  
 TEFKKIKVLGSGAFGTVYKGLWIPEGEKVKIPVAIKELREATSPKANKEILDEAYVMASV  
**>DGM PDGF RECEPTOR, MOUSE KINASE REGION**  
 DQLVLGRTLGSAGFGQVVEATAHGLSHSQATMKVAVKMLKSTARSSEKQALMSELYGDLV  
**>FES THIS IS VFES TYROSINE KINASE**  
 VLNRAVPKDKWVLNHEDLVLGEQIGRGNFGEVFSGRLRADNTLVAVKSCRETLPDIKAK  
**>AF1 RAF1 HUMAN C-RAF-1 ONCOGENE**  
 SEVMLSTRIGSGSFGTVYKKGKWHGDVAVKI LKVVDPTPEQFQAFRNEVAVLRKTRHVNI  
**>MOS CMOS HUMAN C-MOS ONCOGENE**  
 EQVCLLQRLGAGGFGSVYKATYRGVPVAIKQVNKCTKNRLASRRSFWAELNVARLRHDNI  
**>SVK HSVK HERPES SIMPLEX VIRUS PUTATIVE PROTEIN KINASE**  
 MGFTIHGALTPGSEGCVFDSSHPDYPQRVIVKAGWYTSTSHEARLLRRLDHPAILPLLDL

# Multiple Alignment created using ClustalW; Colors Added using BoxShade



BLOCKS Server located blocks

Taking this alignment, blocks can be generated using the BLOCKS server:

```
ID    x6676xbli; BLOCK
AC    x6676xbliA; distance from previous blocks=(1,1)
DE    ../tmp/6676.blin
BL    UNK motif; width=24; seqs=10; 99.5%=0; strength=0
AF1   ( 1) SEVMLSTRIGSGSFGTVYK GKWHG 41
MOS   ( 1) EQVCLLQRLGAGGFSGSVYKATYRG 48
DGM   ( 1) DQLVLGRTLGS GAFGQVVEATAHG 49
GFR   ( 1) TEFKKIKVLGSGAFGTVYKGLWIP 41
D28   ( 1) ANYKRLEKVGEGTYGVVYKALDLR 61
SKH   ( 1) AKYDIKALIGRGSF SRVVRVEHRA 54
APK   ( 1) DQFERIKTLGTGSFGRV MLVKHME 46
EE1   ( 1) TRFRNV TLLGSGEFSEVFQVEDPV 55
FES   ( 1) LNRAVPKDKWVLN HEDLVLGEQIG 100
SVK   ( 1) MGFTIHGALTPGSEGC VFDSSHPD 73
//
```

Taking this alignment, blocks can be generated using the BLOCKS server:

```
ID    x6676xbli; BLOCK
AC    x6676xbliB; distance from previous blocks=(2,2)
DE    ../tmp/6676.blin
BL    UNK motif; width=28; seqs=10; 99.5%=0; strength=0
AF1   ( 27) AVKILKVVDPTPEQFQAFRNEVAVLRKT 87
MOS   ( 27) PVAIKQVNKCTKNRLASRRSFWAELNVA 75
DGM   ( 27) SHSQATMKVAVKMLKSTARSSEKQALMS 92
GFR   ( 27) GEKVKIPVAIKELREATSPKANKEILDE 83
D28   ( 27) PGQGQRVVALKKIRLESEDEGVPSTAIR 83
SKH   ( 27) RQPYAIKMIETKYREGREVCESLRVLR 74
APK   ( 27) GNHYAMKILDKQKVVKLKQIEHTLNEKR 85
EE1   ( 27) TLKYAVKKLKVKFSGPKERNRLLQEVSI 77
FES   ( 27) GNFGEVFSGRLLRADNTLVAVKSCRETLP 100
SVK   ( 27) PQRVIVKAGWYTSTSHEARLLRRLDHPA 92
//
```

# Statistical Methods

- Commonly used methods for locating motifs:
  - Expectation-Maximization (EM)
  - Gibbs Sampling

# Expectation-Maximization

- In the expectation-maximization algorithms,
  - the starting point is a set of sequences expected to have a common sequence pattern that may not be easily detectible.
  - An initial guess is made as to the location and size of the site of interest in each of the sequences.
  - These initial sites are then aligned.
  - Approximate length of signal must be given
- Randomly assign locations of this motif in each sequence

# Expectation-Maximization

- Two steps:
  - Expectation Step
  - Maximization Step

# Expectation-Maximization

- Expectation step
  - In the expectation step, background residue frequencies are calculated based on those residues that are not in the initially aligned sites.
  - Column specific residues are calculated for each position in the initial motif alignment.
  - Using this information, the probability of finding the site at any position in the sequences can then be calculated.
  - Residues not in a motif are background
- Frequencies used to determine probability of finding site at any position in a sequence to fit motif model

# Maximization Step

- Maximization step
  - In the maximization step, the counts of residues for each position in the site as found in the expectation step are used to calculate the location within each sequence that maximally aligns to the motif pattern calculated in the expectation step.
  - This is done for each of the sequences.
  - Once a new motif location has been calculated, the expectation step is repeated.
  - This cycle continues until the solution converges.

TCAGAACCAGTTATAA**ATTTAT**CATTTCCTTCTCCACTCCT  
CCCACGCA**GCCGCC**CTCTCCCCGGTCACTGACTGGTCCTG  
TCGACCTCTGAACCTATCAGGGACCA**CAGTCA**GCCAGGCAAG  
AAAACACTTGAG**GGAGCA**GATAACTGGGCCAACCATGACTC  
GGGTGAATGGTACTGCT**GATTACA**ACCTCTGGTGCTGC  
AGCCTAGAGT**GATGAC**TCCTATCTGGGTCCCCAGCAGGA  
GCCTCAGGATCCAGCACACAT**TATCAC**AAACTTAGTGTCCA  
CATTATCAC**AAACTT**AGTGTCCATCCATCACTGCTGACCCCT  
TCGGAACAAGGCAAA**GGCTAT**AAAAAAATTAAGCAGC  
GCCCCTCCCCA**CACTAT**CTCAATGCAAATATCTGTCTGAAACGGTTCC  
CATGCCCTCAAGTGTGCAGATTGGT**CACAGC**ATTTCAAGG  
GATTGGTCACAGCAT**TTCAAG**GGAGAGACCTCATTGTAAG  
TCCCCAACTCCCAACTGACCTTAT**CTGTGG**GGGAGGCTTTTGA  
CCTTATCTGT**GGGGGA**GGCTTTTGAAGTAATTAGGTTTAGC  
ATTATTTTCCTTATCAGAAGC**AGAGAG**ACAAGCCATTTCTCTTTCCCTCCCGGT  
AGG**CTATAA**AAAAAATTAAGCAGCAGTATCCTCTTGGGGGCCCTTC  
CCAGCACACACTTATC**CAGTGG**TAAATACACATCAT  
TCAAATAGGTACGGATAAG**TAGATA**TTGAAGTAAGGAT  
ACTTGGGGTTCCAGTTTGATAAGAAAAGACTT**CCTGTG**GGA  
TGGCCGC**AGGAAG**GTGGCCTGGAAGATAACAGCTAGTAGGCTAAGGCCAG  
**CAACCA**CAACCTCTGTATCCGGTAGTGGCAGATGGAAA  
CTGTATCCGGTAG**TGGCAG**ATGGAAGAGAAAACGGTTAGAA  
GAAAAAAATAAATGAAGTCTGCC**TATCTC**CGGGCCAGAGCCCT  
TGCCTTGTCTGTTGTAGATAATGAATCTATCCTCCA**GTGACT**  
GGCCAGGCTGAT**GGGCCT**TATCTCTTACCACCTGGCTGT  
CAACAGCAGGTCTACTATCGCCTCCCTCT**AGTCTC**TG  
CCAACCG**TTAATG**CTAGAGTTATCACTTTCTGTTATCAAGTGGCTTCAGCTATGCA  
GGGAGGGTGGGGCCCTATCTCTCTA**GACTCT**GTG  
CTTTGT**ACTGGA**TCTGATAAGAAACACCACCCCTGC

Example of EM:  
begin with an  
initial, Random  
alignment:

# Residue Counts

- From this alignment, the frequency of each base occurring is calculated.
- In this case, the motif we are searching for is six bases wide.
- Therefore, we need to calculate seven different sets of frequencies:
  - One for the background, and
  - one for each of the columns in the motif.
- Calculating the total counts, we get:

Nucleotide	Motif Position (0 = Background)						
	0	1	2	3	4	5	6
A	279	6	12	6	6	11	7
C	280	8	3	5	7	7	7
G	225	9	8	10	7	5	8
T	262	6	6	8	9	6	7

**Table 1: Calculation of observed counts for initial alignment of figure 1**

# Residue Frequencies

- After calculating the observed counts for each of the positions, we can convert these to observed frequencies:

Nucleotide	Motif Position (0 = Background)						
	0	1	2	3	4	5	6
<b>A</b>	.267	.256	.296	.256	.256	.289	.263
<b>C</b>	.267	.263	.230	.243	.256	.256	.256
<b>G</b>	.216	.240	.233	.246	.226	.213	.233
<b>T</b>	.250	.241	.241	.254	.261	.241	.248

**Table 2: Calculation of residue frequencies for initial alignment of figure 1**

# Example Maximization Step

- In the expectation step, the residue frequencies for the motif are used to estimate the composition of the motif site.
- The expectation step attempts to maximally discriminate between sequence within and not within the site.
- For each sequence, each possible motif location is considered in order to find the most probable location given the current motif.
- Consider the first sequence:
- TCAGAACCAGTTATAA**ATTTAT**CATTCCTTCTCCACTCCT
- There are 41 residues;  $41-6+1 = 36$  sites to consider

	1	2	3	4	5	6	1*2*3*4*5* 6	RANDOM	ODDS
TCAGAA	.241	.230	.256	.226	.289	.263	0.000244	0.000274	0.89
CAGAAC	.263	.296	.246	.256	.289	.256	0.000363	0.000362	1.00
AGAACC	.256	.233	.256	.256	.256	.256	0.000256	0.000362	0.71
GAACCA	.240	.296	.256	.256	.256	.263	0.000313	0.000362	0.87
AACCAG	.256	.296	.243	.256	.289	.233	0.000317	0.000362	0.88
ACCAGT	.256	.230	.243	.256	.213	.248	0.000193	0.000274	0.71
CCAGTT	.263	.230	.256	.226	.241	.248	0.000209	0.000257	0.81
<b>CAGTTA</b>	<b>.263</b>	<b>.296</b>	<b>.246</b>	<b>.261</b>	<b>.241</b>	<b>.263</b>	<b>0.000317</b>	<b>0.000257</b>	<b>1.23</b>
AGTTAT	.256	.233	.254	.261	.289	.248	0.000283	0.000241	1.18
GTTATA	.240	.241	.254	.256	.241	.263	0.000238	0.000241	0.99
TTATAA	.241	.241	.256	.261	.289	.263	0.000295	0.000297	0.99
TATAAA	.241	.296	.254	.256	.289	.263	0.000353	0.000297	1.19
ATAAAT	.256	.241	.256	.256	.289	.248	0.000290	0.000318	0.91
TAAATT	.241	.296	.256	.256	.241	.248	0.000279	0.000297	0.94
AAATTT	.256	.296	.256	.261	.241	.248	0.000303	0.000297	1.02
AATTTA	.256	.296	.254	.261	.241	.263	0.000318	0.000297	1.07
ATTTAT	.256	.241	.254	.261	.289	.248	0.000293	0.000278	1.05
TTTATC	.241	.241	.254	.256	.241	.256	0.000233	0.000278	0.84

TTATCA	.241	.241	.256	.261	.256	.263	0.000261	0.000297	0.88
TATCAT	.241	.296	.254	.256	.289	.248	0.000332	0.000297	1.12
ATCATT	.256	.241	.243	.256	.241	.248	0.000229	0.000297	0.77
TCATTT	.241	.230	.256	.261	.241	.248	0.000221	0.000278	0.80
CATTTT	.263	.296	.254	.261	.241	.256	0.000318	0.000297	1.07
ATTTCC	.256	.241	.254	.261	.256	.256	0.000268	0.000297	0.90
TTCCTT	.241	.241	.254	.256	.256	.248	0.000240	0.000278	0.86
TTCCTT	.241	.241	.243	.256	.241	.248	0.000216	0.000278	0.78
TCCTTC	.241	.230	.243	.261	.241	.256	0.000217	0.000297	0.73
CCTTCT	.263	.230	.254	.261	.256	.248	0.000255	0.000297	0.86
CTTCTC	.263	.241	.254	.256	.241	.256	0.000254	0.000297	0.86
TTCTCC	.241	.241	.243	.261	.256	.256	0.000241	0.000297	0.81
TCTCCA	.241	.230	.254	.256	.256	.263	0.000243	0.000318	0.76
CTCCAC	.263	.241	.243	.256	.289	.256	0.000292	0.000339	0.86
TCCACT	.241	.230	.243	.256	.256	.248	0.000219	0.000318	0.69
CCACTC	.263	.230	.256	.256	.241	.256	0.000245	0.000339	0.72
CACTCC	.263	.296	.243	.261	.256	.256	0.000324	0.000339	0.95
ACTCCT	.256	.230	.254	.256	.256	.248	0.000243	0.000318	0.76

- The six base site **CAGTTA** beginning at base 8 is calculated to have the highest odds probability.
- Therefore, it is chosen as the new site in sequence 1.
- This is repeated for each of the sequences.
- In the maximization step, the newly chosen sites for each of the sequences are used to recalculate the frequency table.
- The expectation/maximization cycle is then repeated, until the results converge on a set of motifs.

# Maximization Step

- Before: Random Alignment
- TCAGAACCAGTTATAA**ATTTAT**CATTTTCCTTCTCCACTCCT
- After: Maximal location (given random motif alignment) (first round)
- TCAGAAC**CAGTTA**TAA**ATTTAT**CATTTTCCTTCTCCACTCCT

# Available E-M Programs

- MEME – Uses E-M algorithms as explained
- **Multiple EM for Motif Elicitation (MEME)** is a program developed that uses the expectation-maximization methods as described previously.
- ParaMEME searches for blocks using the EM algorithm, while MetaMEME searches for profiles using Hidden Markov Models.
- MEME locates one or more ungapped patterns in a single DNA or protein sequence, or in a series of sequences.
- A search is conducted on a variety of motif widths in order to determine the most likely width for the profile.
- This likelihood is based on the log likelihood score calculated after the EM algorithm.

# MEME Software

- One of three types of motif models can be chosen:
  - OOPS: One expected occurrence per sequence
  - ZOOPS: Zero or one expected occurrence per sequence
  - TCM: Any number of occurrences of the motif

# MEME Software

- Various prior knowledge can be added to MEME, including the expected number of motifs, the expected length of the motif, and whether or not the motif is palindromic (only applicable for DNA sequences).
  - Palindromic sequences (DNA)
  - Expected number of motifs
  - Expected length of motifs